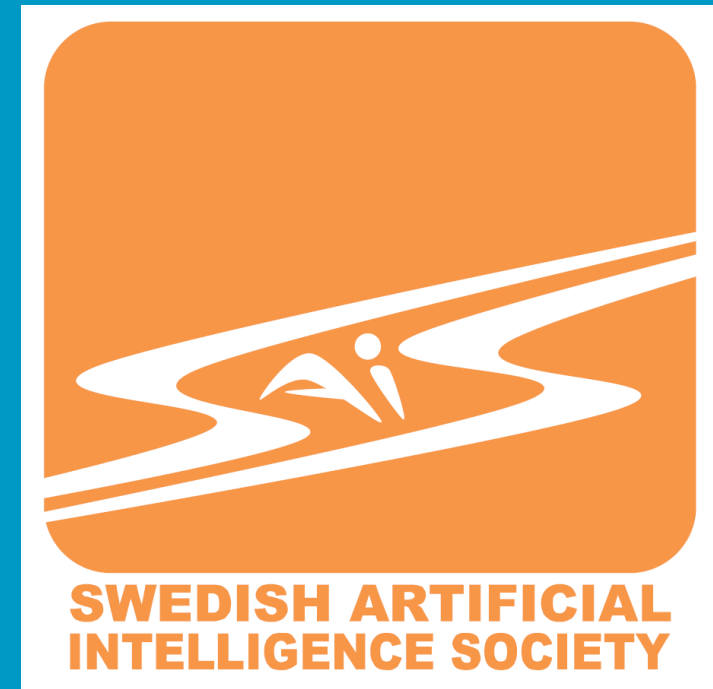
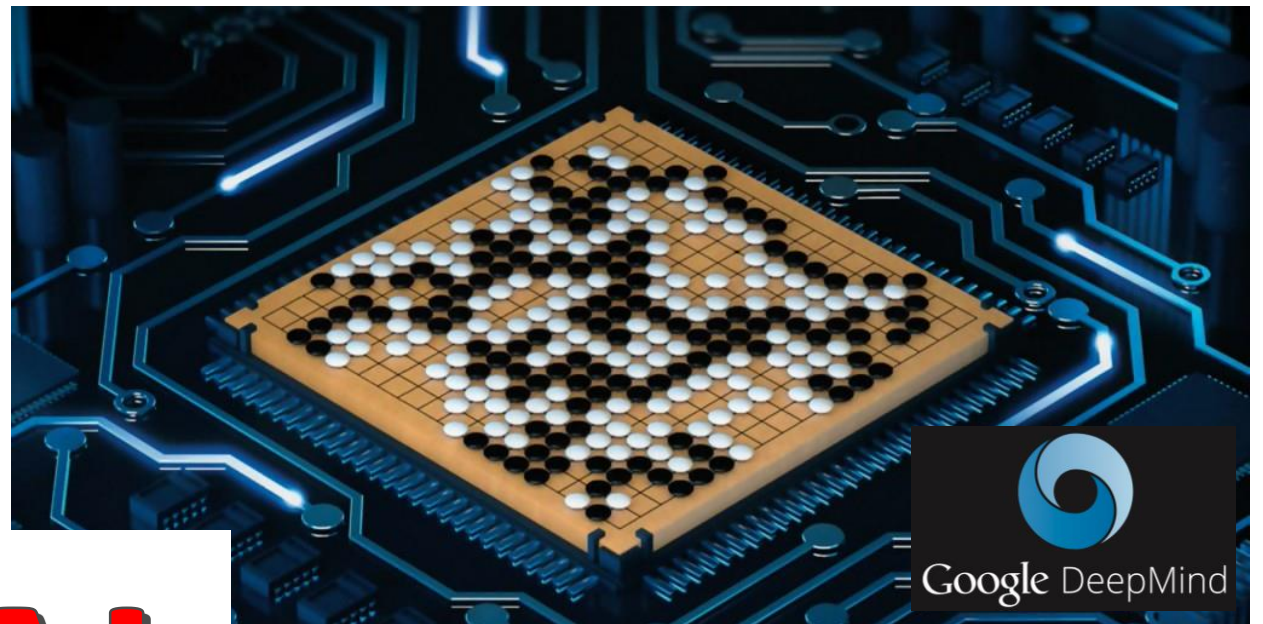


Trustworthy Human-Centric AI

Fredrik Heintz, Dept. of Computer Science
Linköpings universitet
fredrik.heintz@liu.se
[@FredrikHeintz](https://twitter.com/FredrikHeintz)

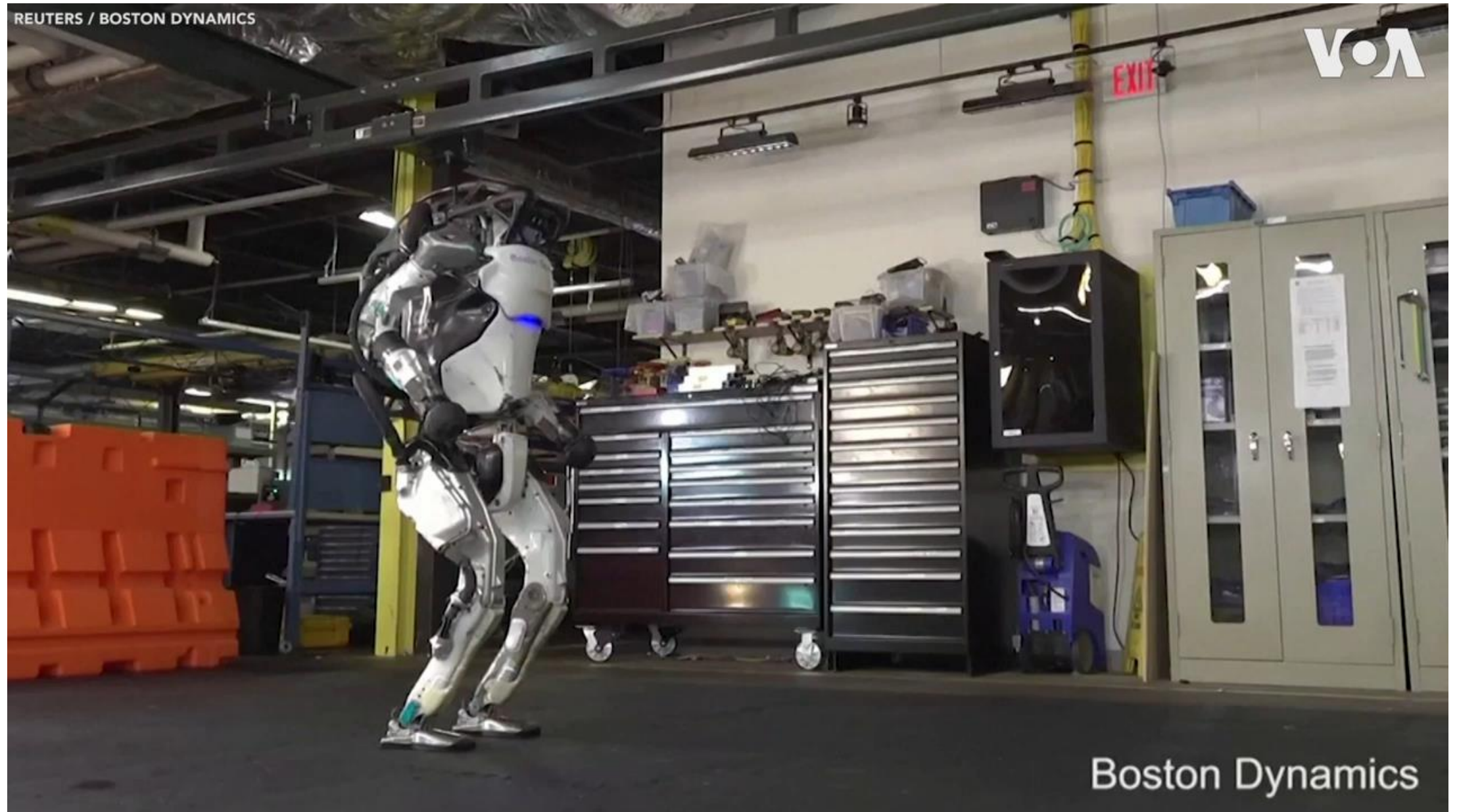




AI







Boston Dynamics





Artificial Intelligence – What is it? – Definitions

“Artificial Intelligence is the **science and engineering of making intelligent machines**, especially intelligent computer programs.”

- John McCarthy, Stanford

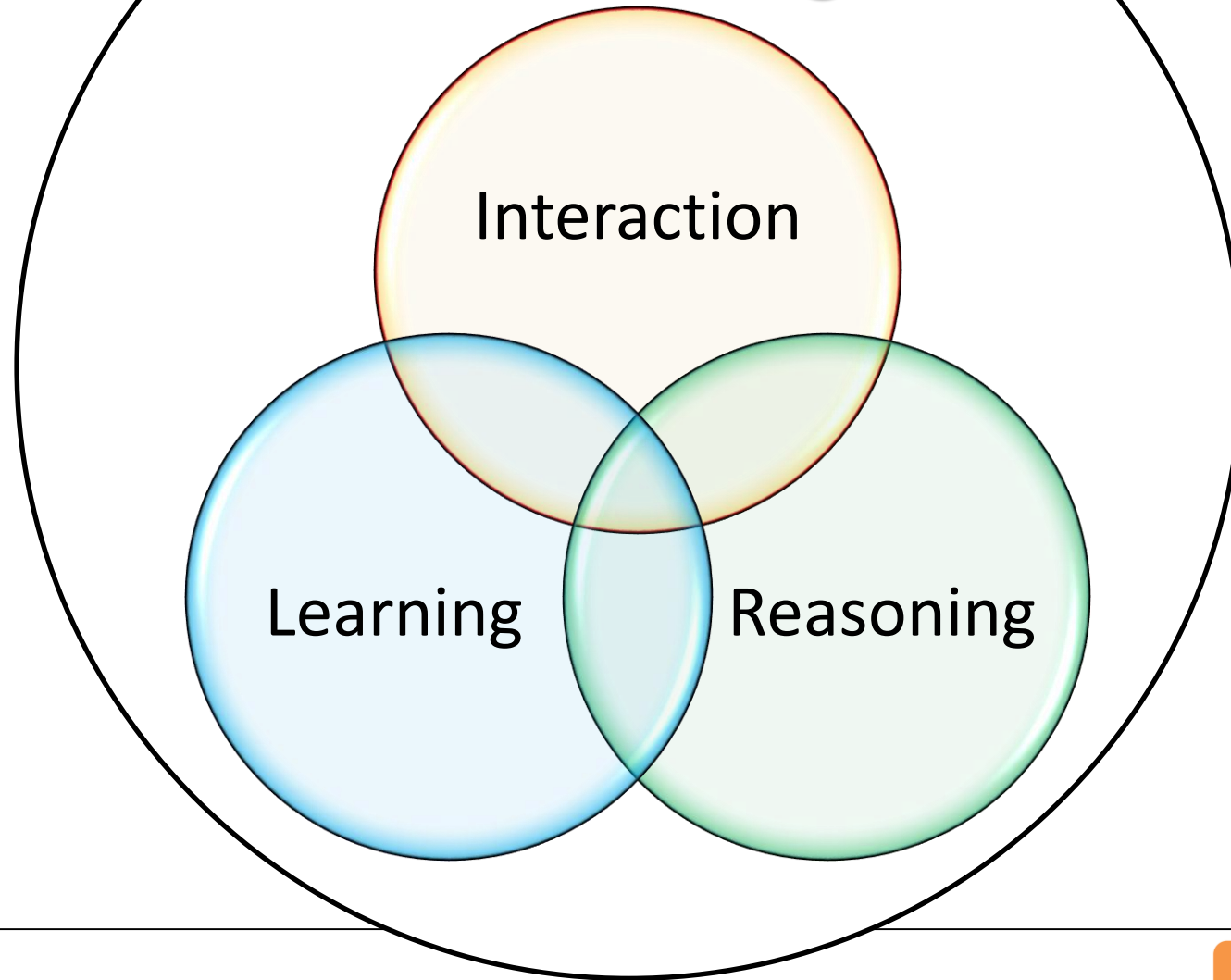
“Artificial intelligence (AI) refers to **systems that display intelligent behaviour** by analysing their environment and taking actions – with some degree of **autonomy** – to achieve specific **goals**.”

- EU Communication 25 April 2018

“the scientific understanding of the **mechanisms underlying thought and intelligent behavior** and their embodiment **in machines**.”

- AAAI

Artificial Intelligence



Digitalization and AI

Digitization

first wave



Big Data

second wave

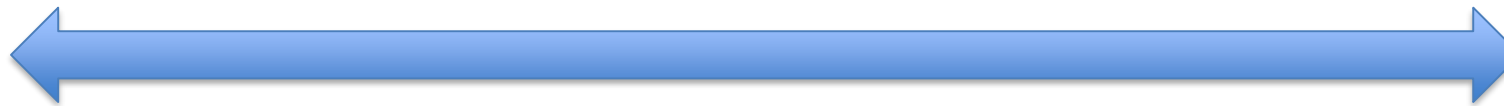


AI

third wave

Digitization

AI



Well defined problems
Predictable situations
Structured data
General solutions
Rationalizes
Evolutionary
...

Hard to define problems
Unanticipated situations
Unstructured data
Adaptable solutions
Amplifies
Revolutionary
...

Skolans Digitalisering

- A. *Skolans digitala infrastruktur*, tillgång till datorer, fungerande nätverk, hantering av inköp av program, supportpersonal mm.
- B. Innehållet i undervisningen för att möta de behov som dagens digitaliserade samhälle kräver, ***vad*** ska undervisas.
- C. Digitalisering av skolans undervisning och pedagogik utifrån de digitala verktyg och möjligheter som finns idag, ***hur*** undervisningen genomförs.
- D. *Digitaliseringen av skolans stödsystem* och övrig verksamhet, t.ex. administrativa system, interaktion med föräldrar, intranät, osv.

Hur kan man dra nytta av AI?

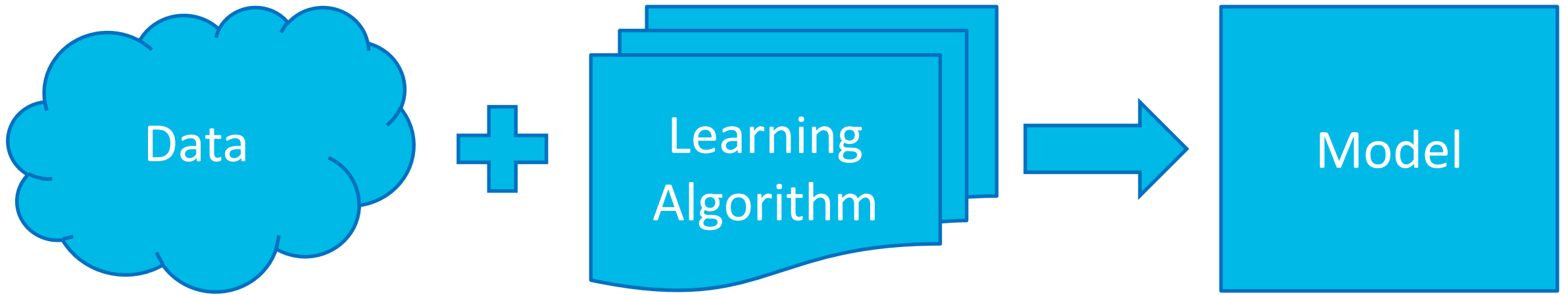
- Textbaserad AI – hantera stora mängder text och dokument
 - Sammanfatta, organisera, analysera text
 - Arbetsförmedlingens nya jobbmatchning
 - Nyhetsvärdering från diariet
- Dialogbaserad AI – hantera kommunikation
 - En ingång till kommunen/verksamheten
 - Bara fråga om information en gång
- Analytisk AI – förstå och förutspå vad som händer
 - Förstå faktisk användning av resurser
 - Förutspå kommande behov

Artificial Intelligence
(AI)

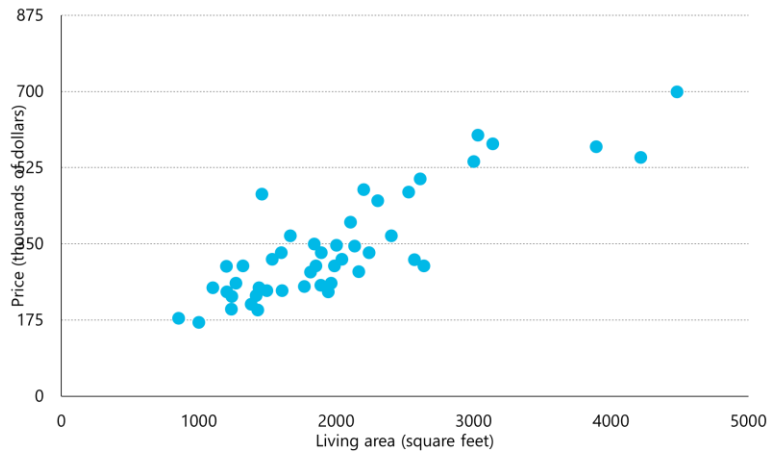
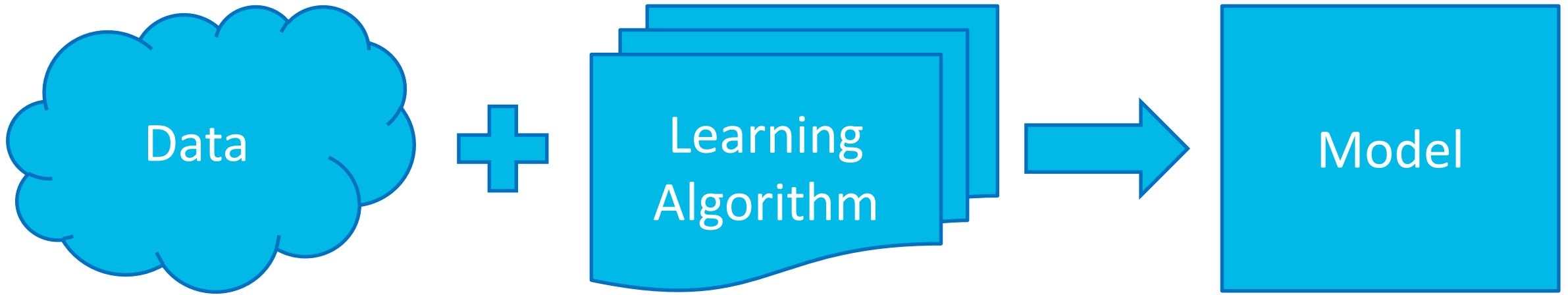
Machine Learning
(ML)

Deep Learning (DL)

Machine Learning

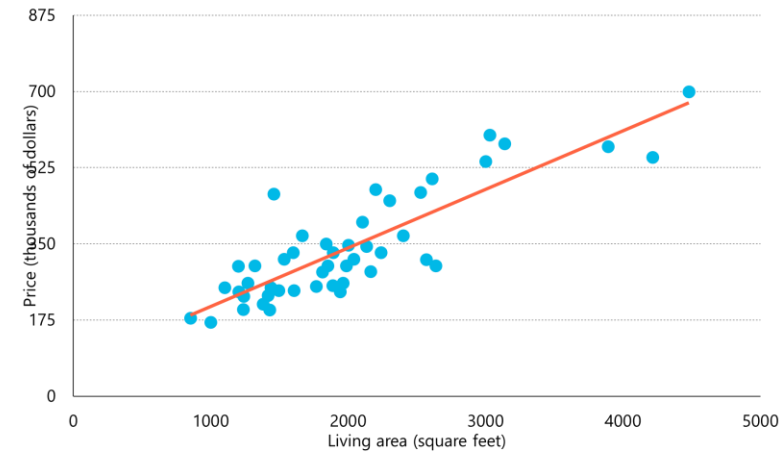


Machine Learning

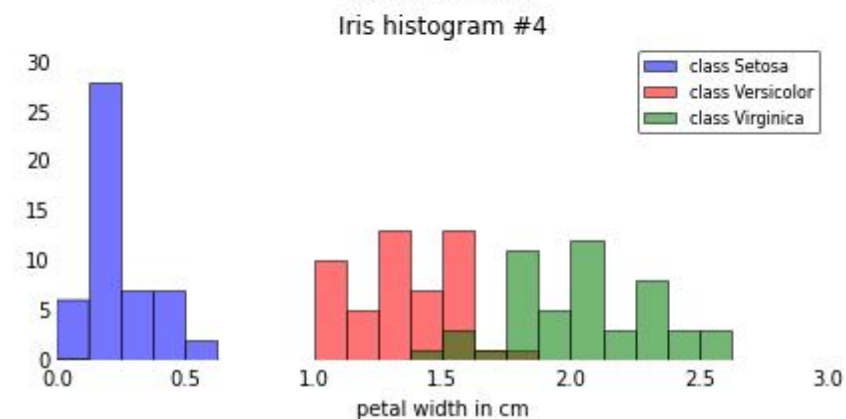
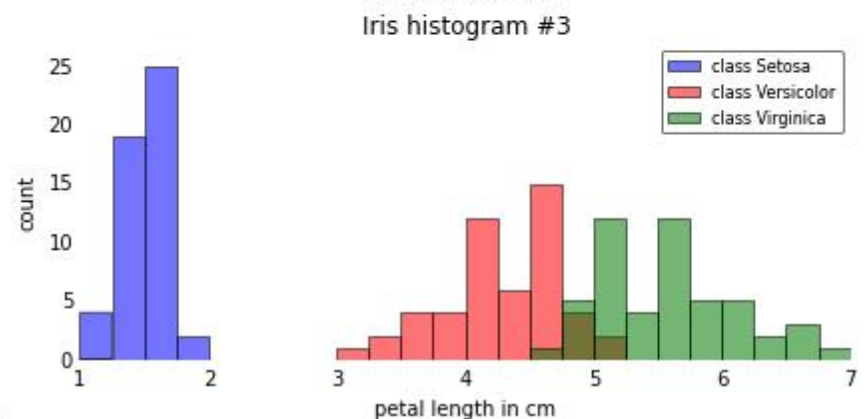
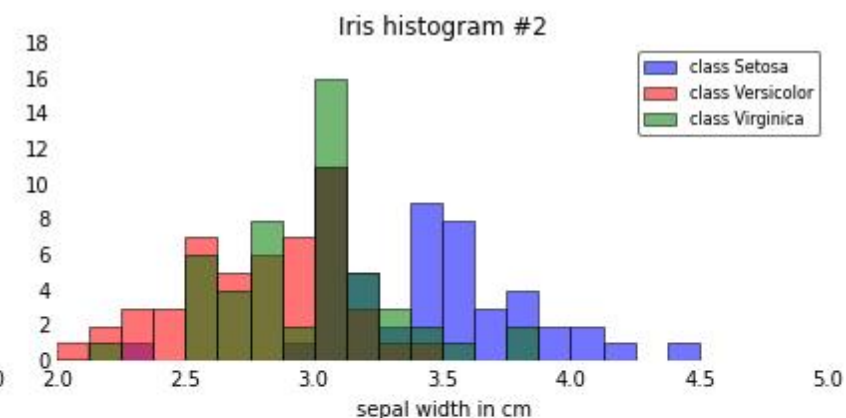
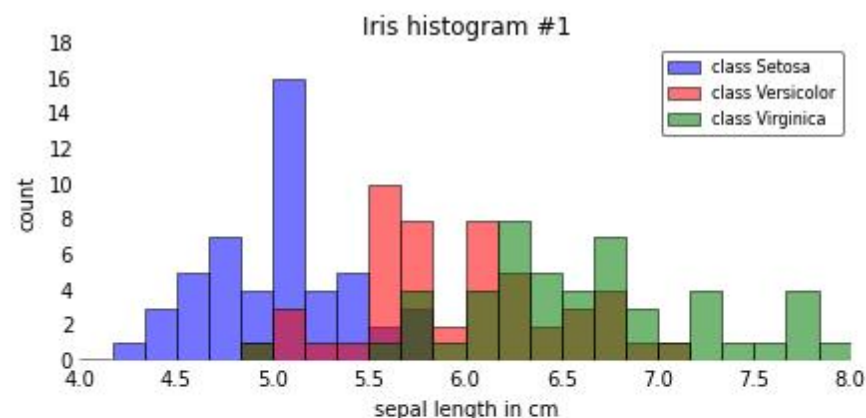
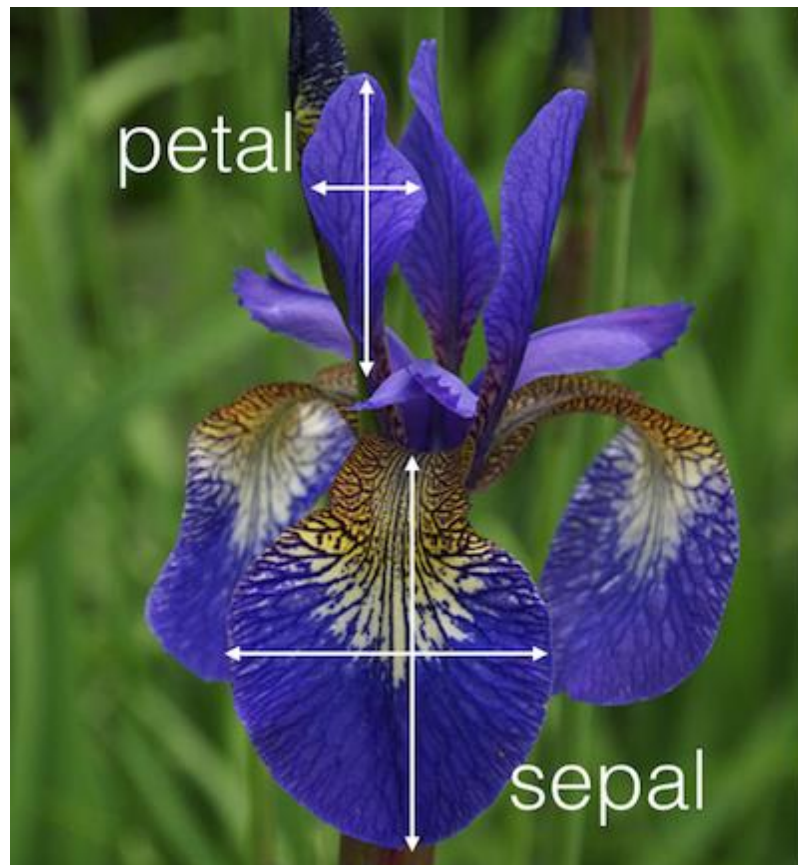


Linear regression

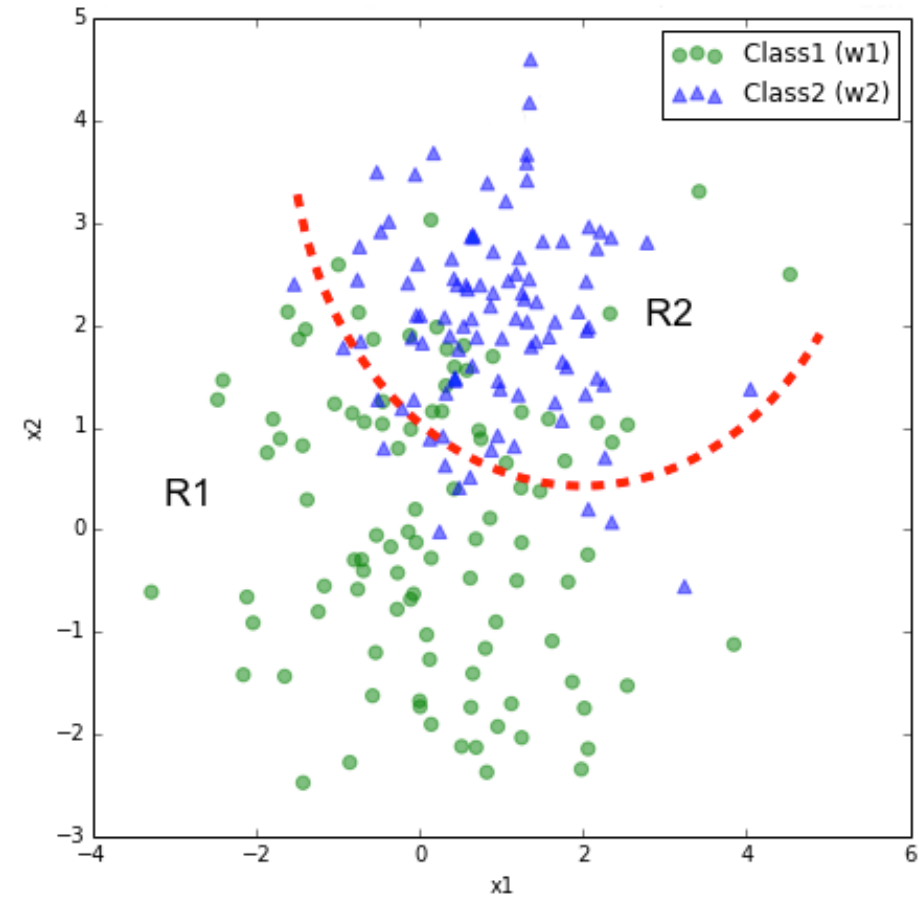
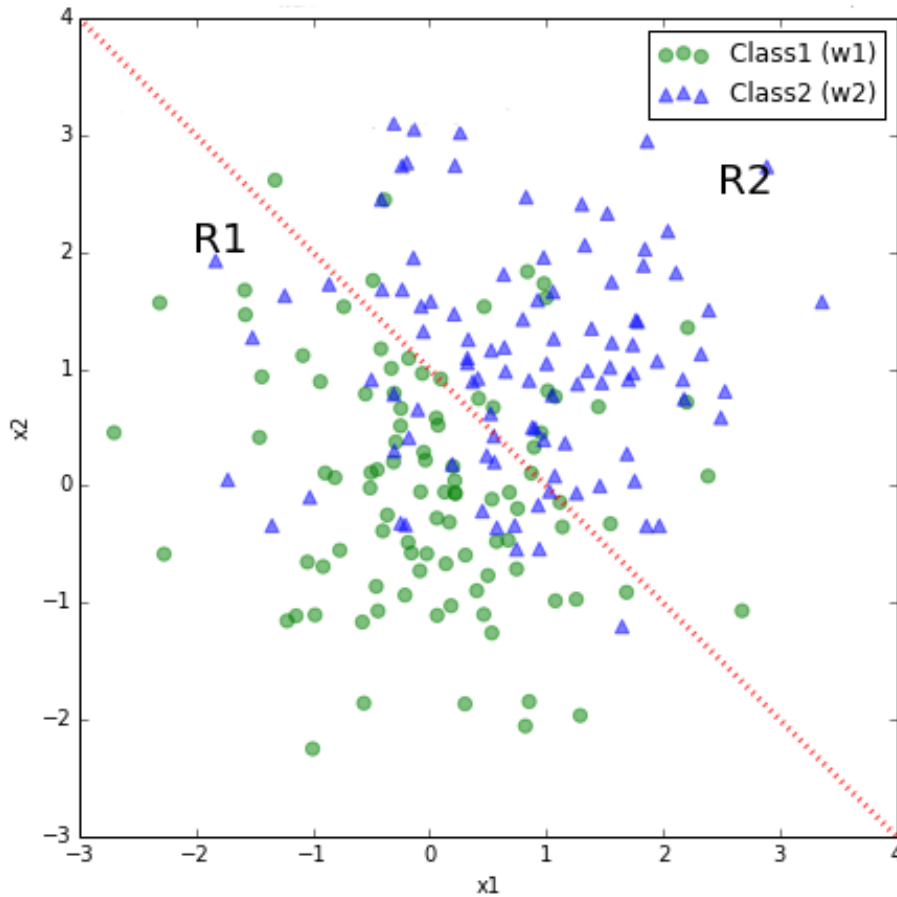
$$\hat{y} = x\theta$$



The Importance of Feature Selection

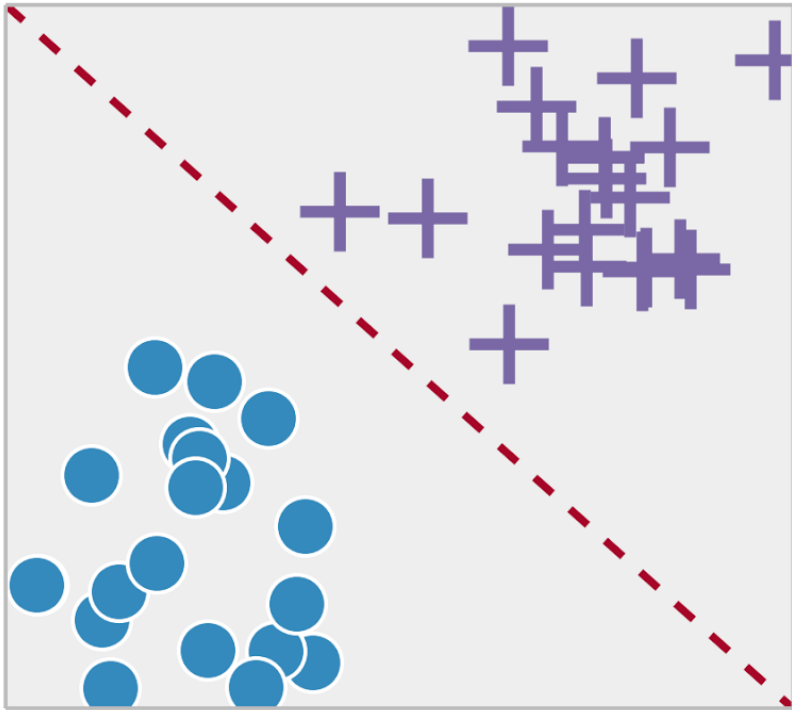


Classification

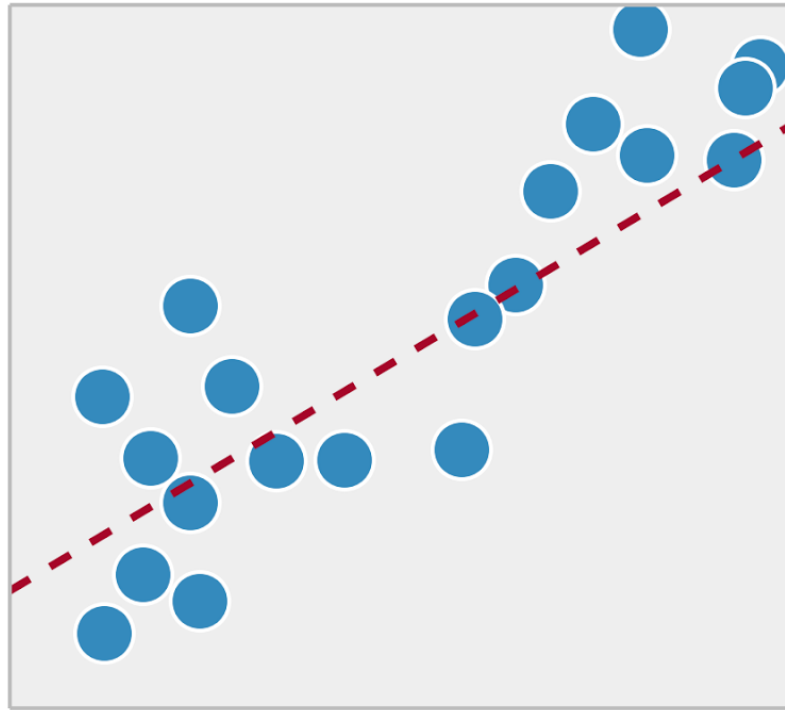


Model Types

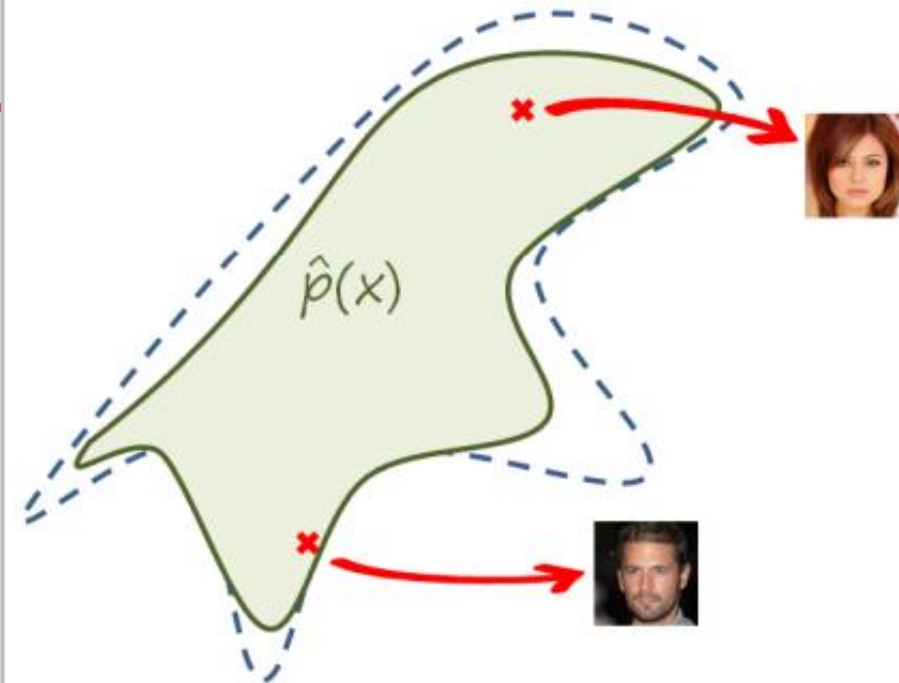
Classification



Regression

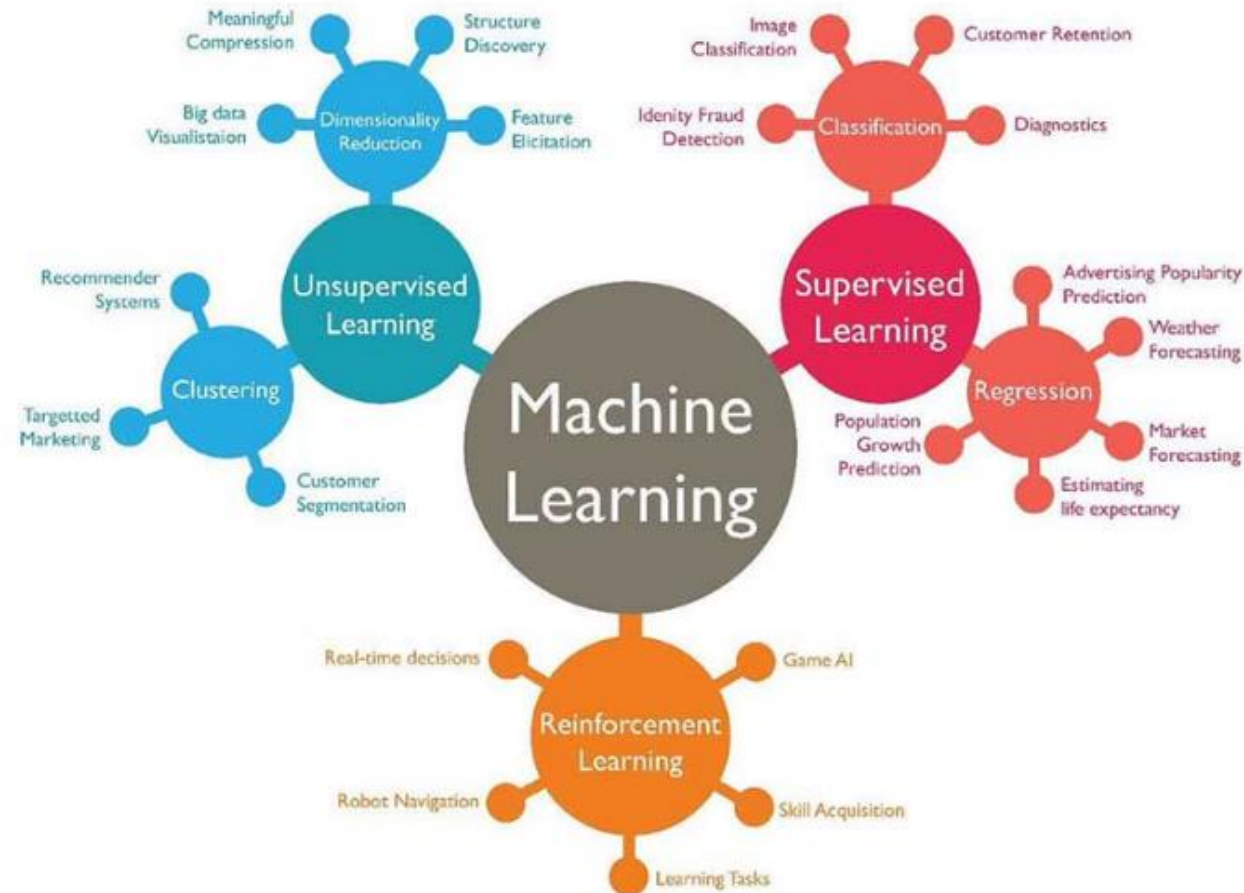


Generative

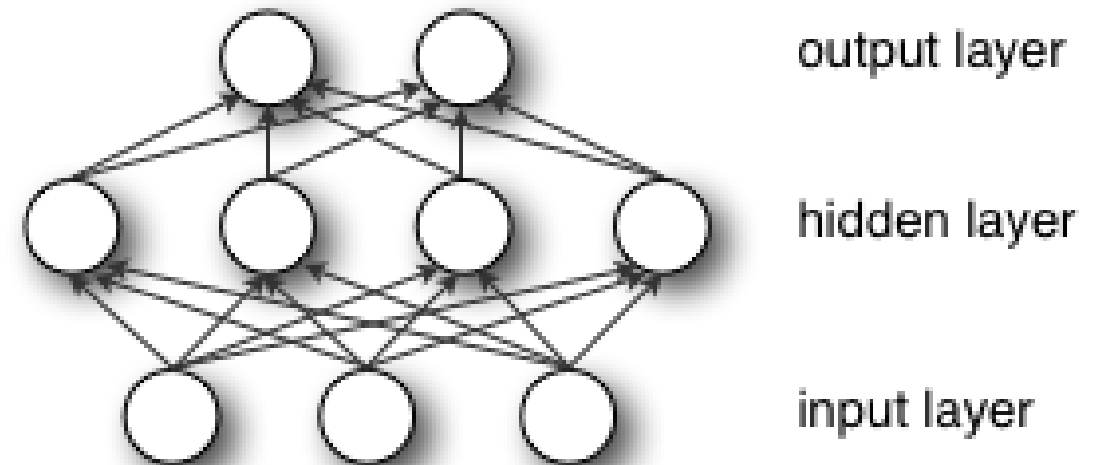
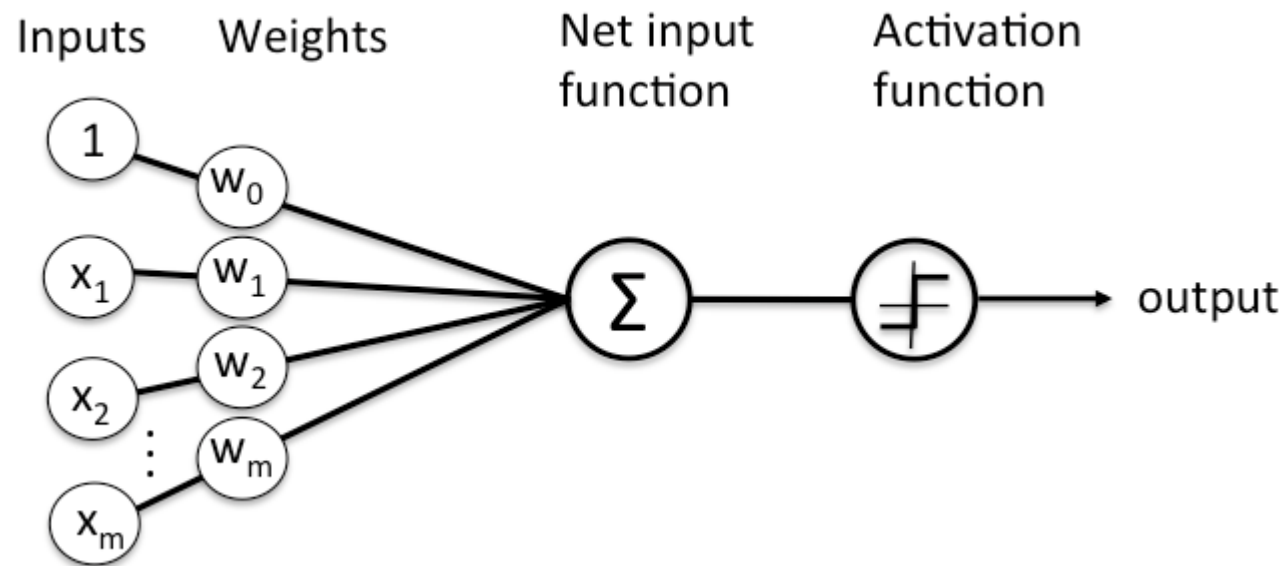


Types of Machine Learning

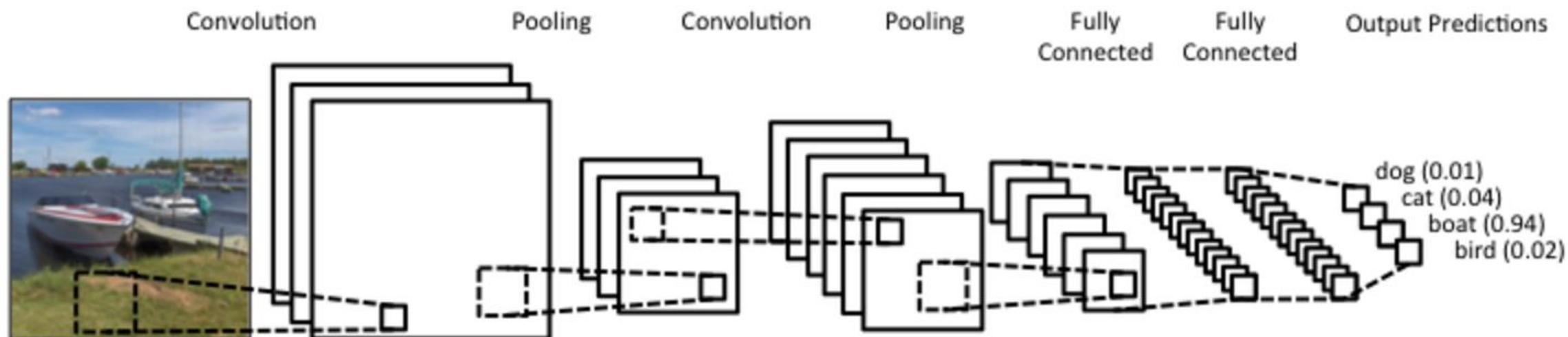
- Supervised learning
 - Given input-output examples $f(X)=Y$, learn the function $f()$.
- Unsupervised learning
 - Given input examples, find patterns such as clusters
- Reinforcement learning
 - Select and execute an action, get feedback, update policy (what action to do in which state).



Neural Networks



Deep Neural Networks



Neural Networks Timeline

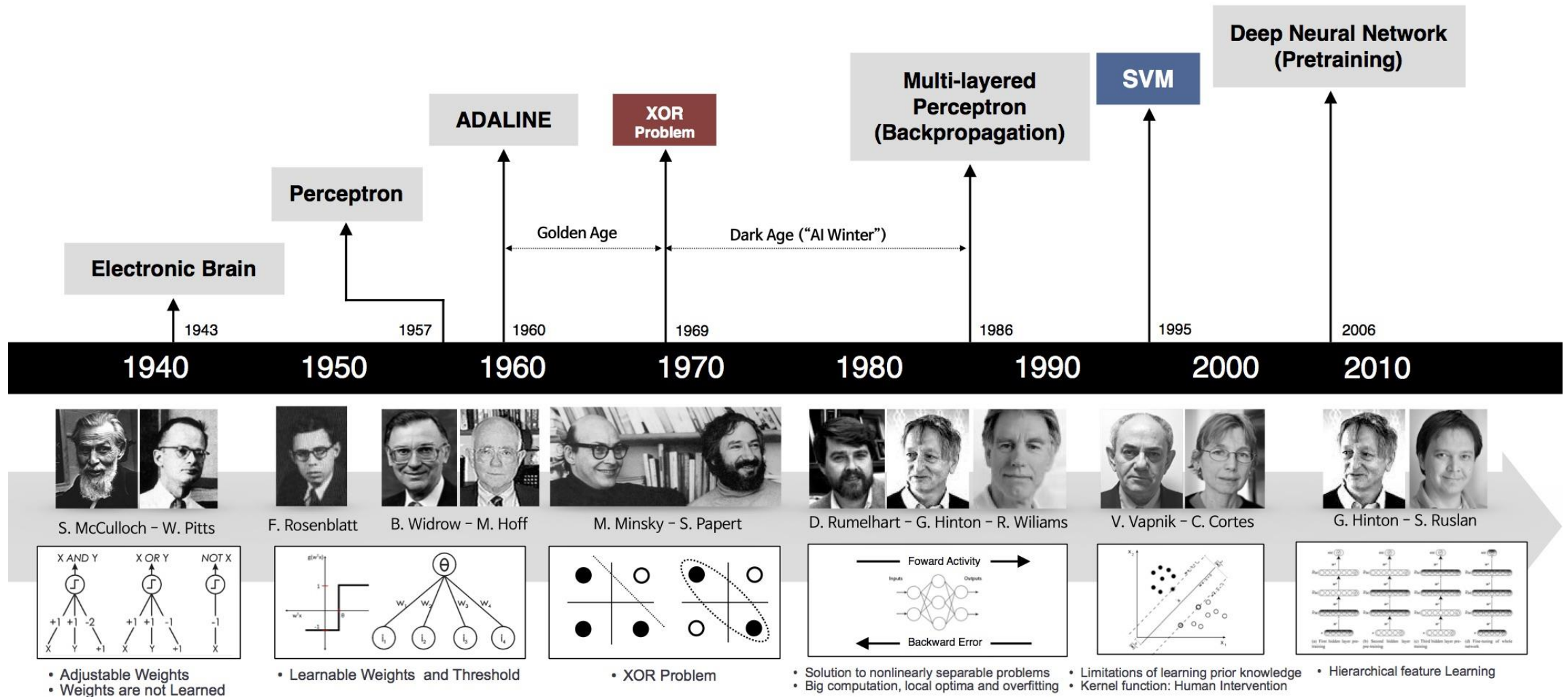
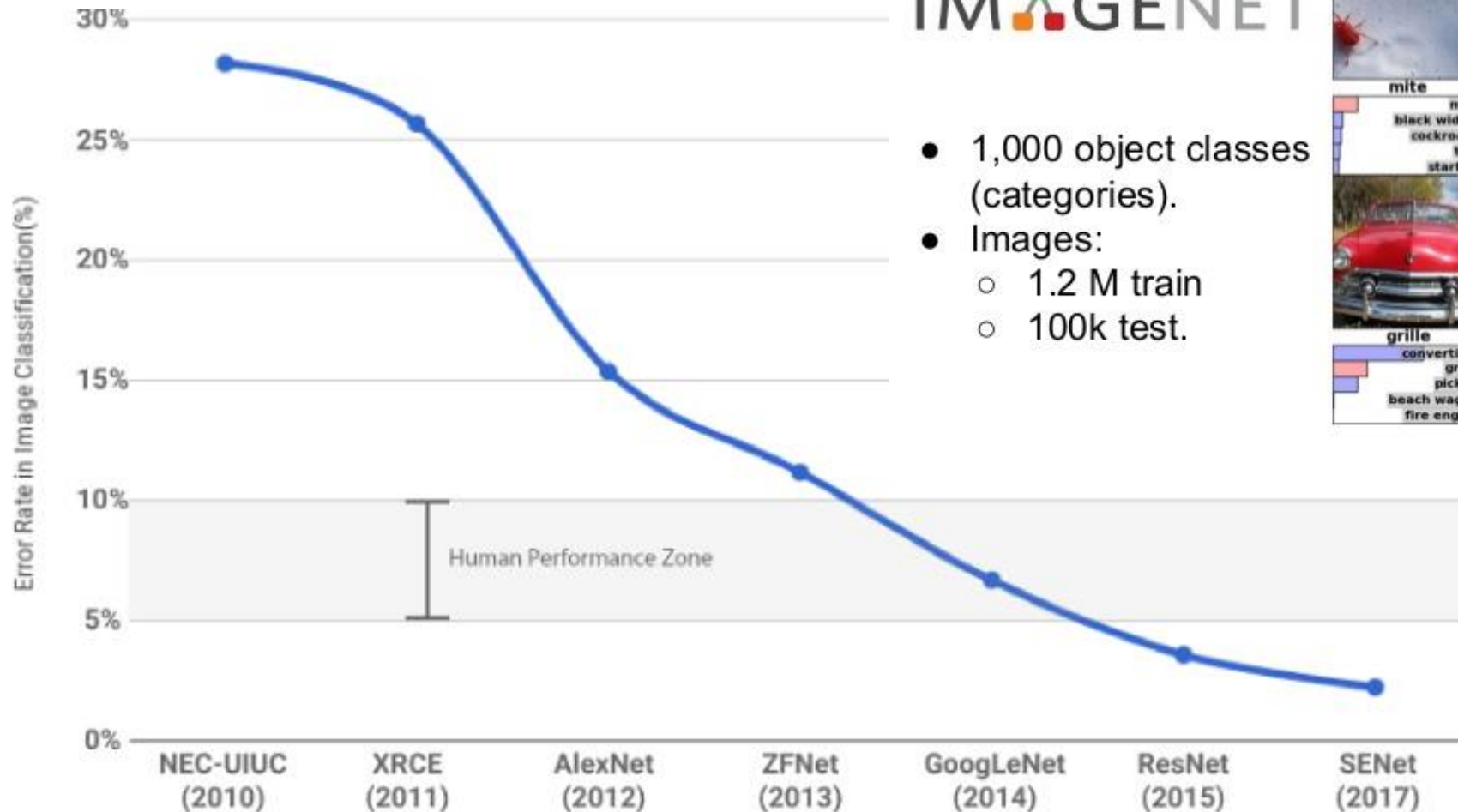


Image Classification



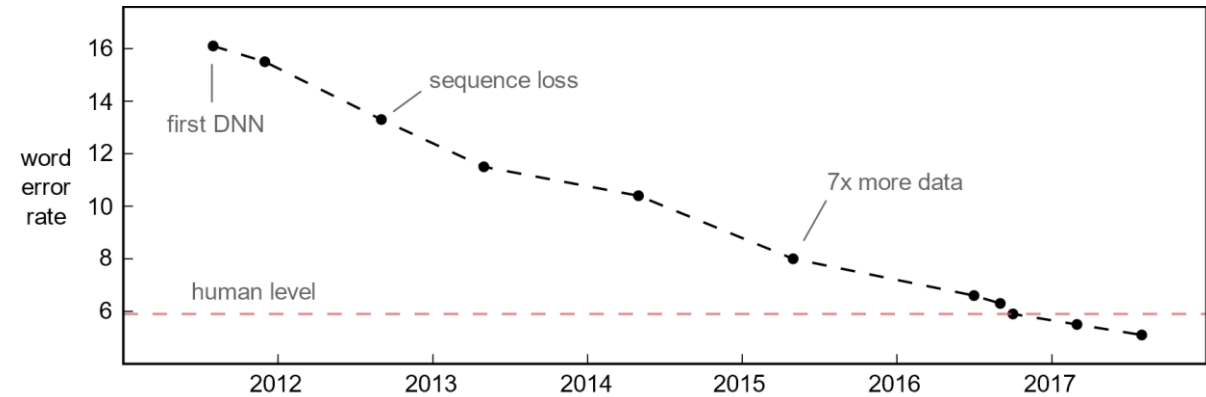
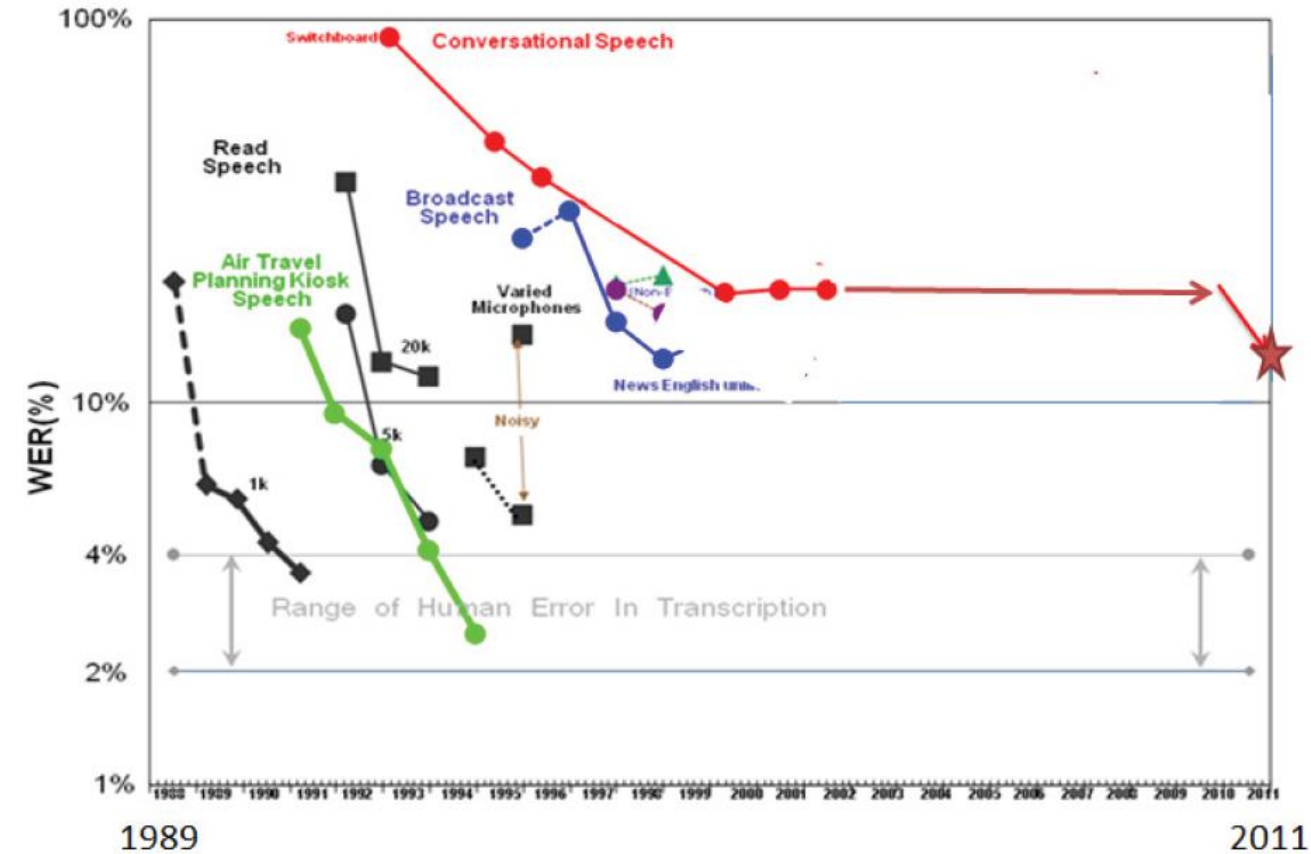
ImageNet Challenge

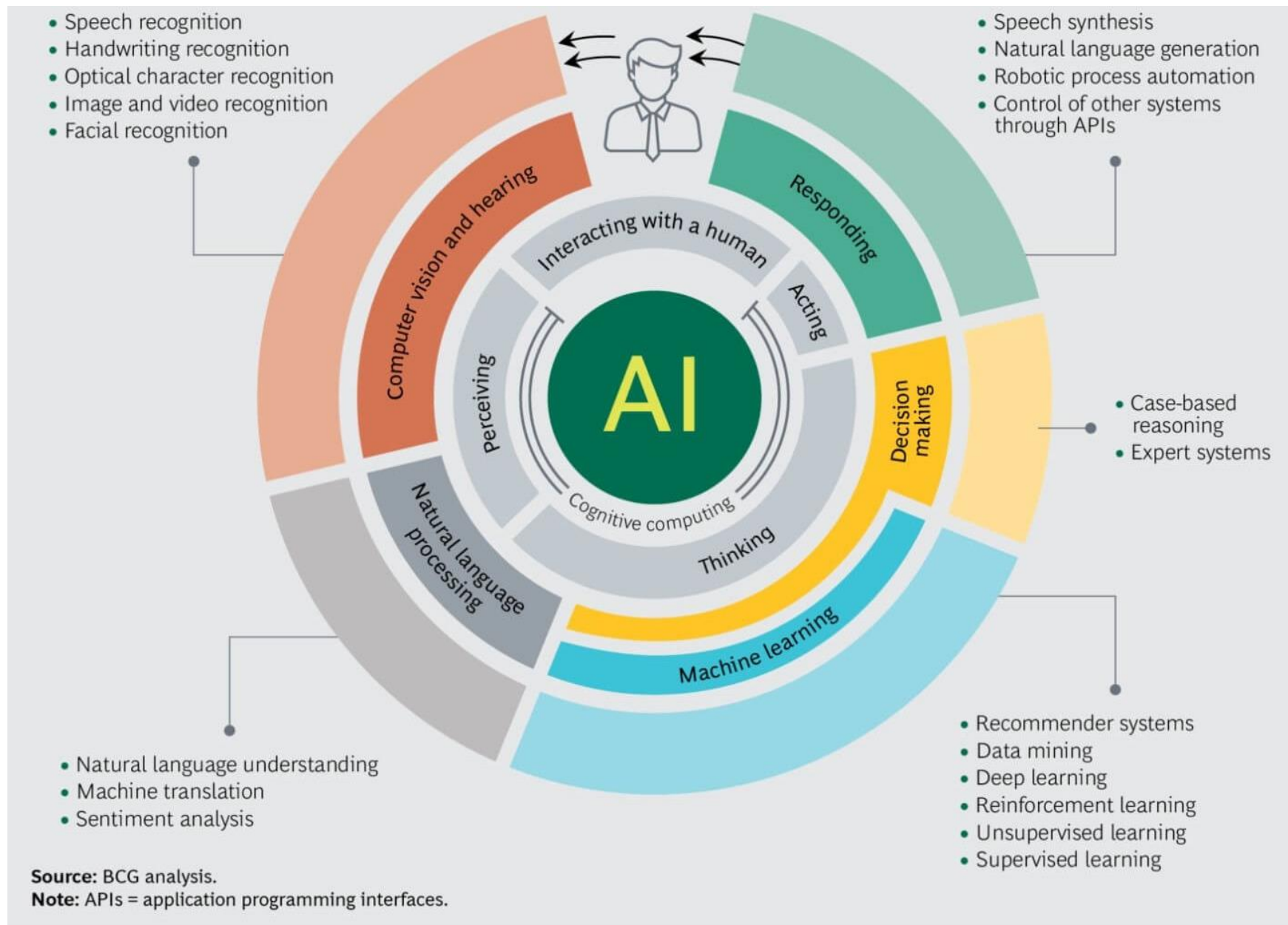
IMAGENET

- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test.

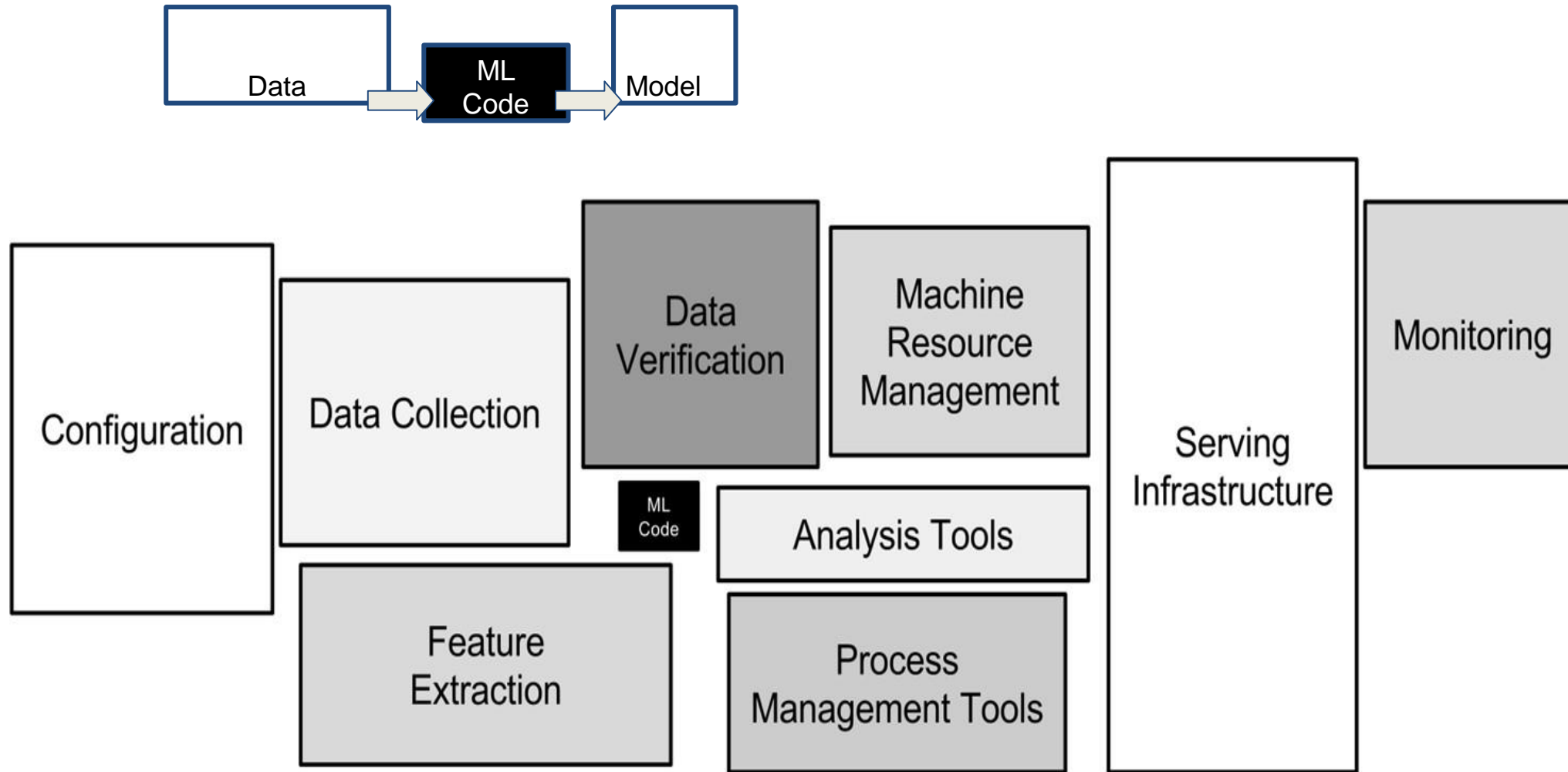


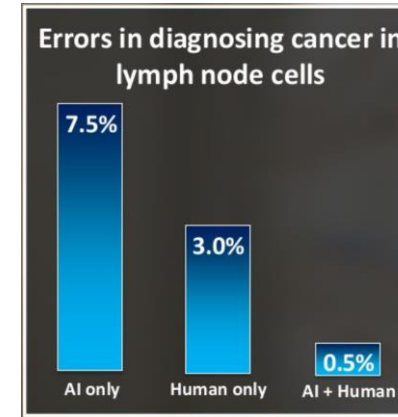
Speech Recognition





The bigger system / picture



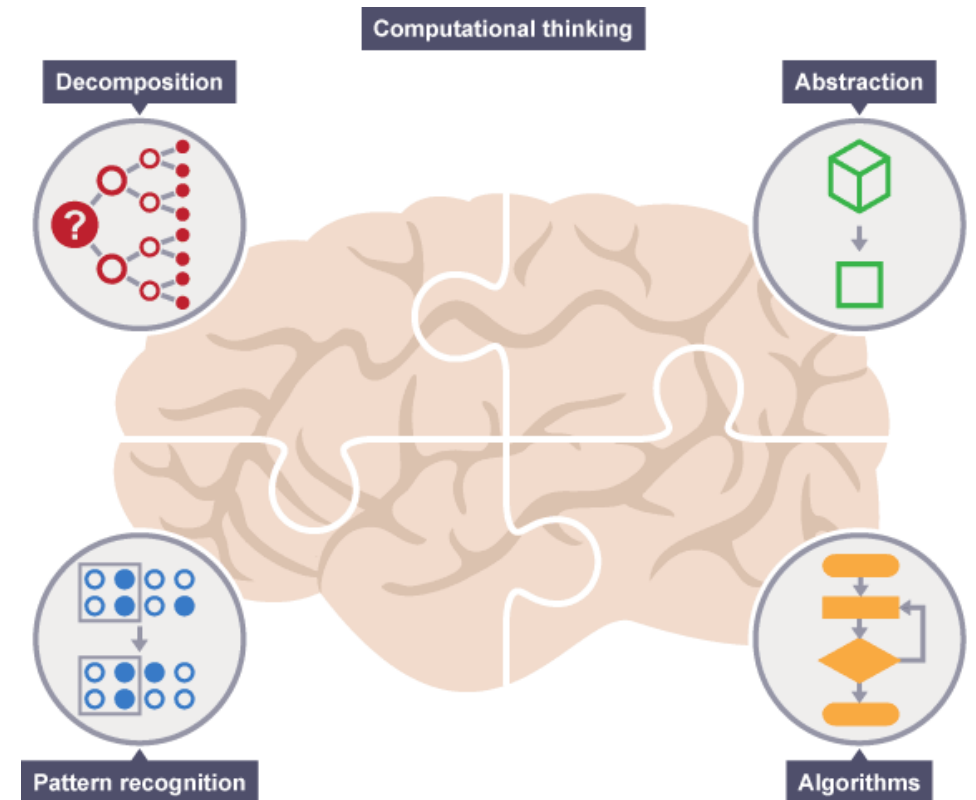


“Weak human + machine + superior process was greater than a strong computer and, remarkably, greater than a strong human + machine with inferior process.”

Garry Kasparov

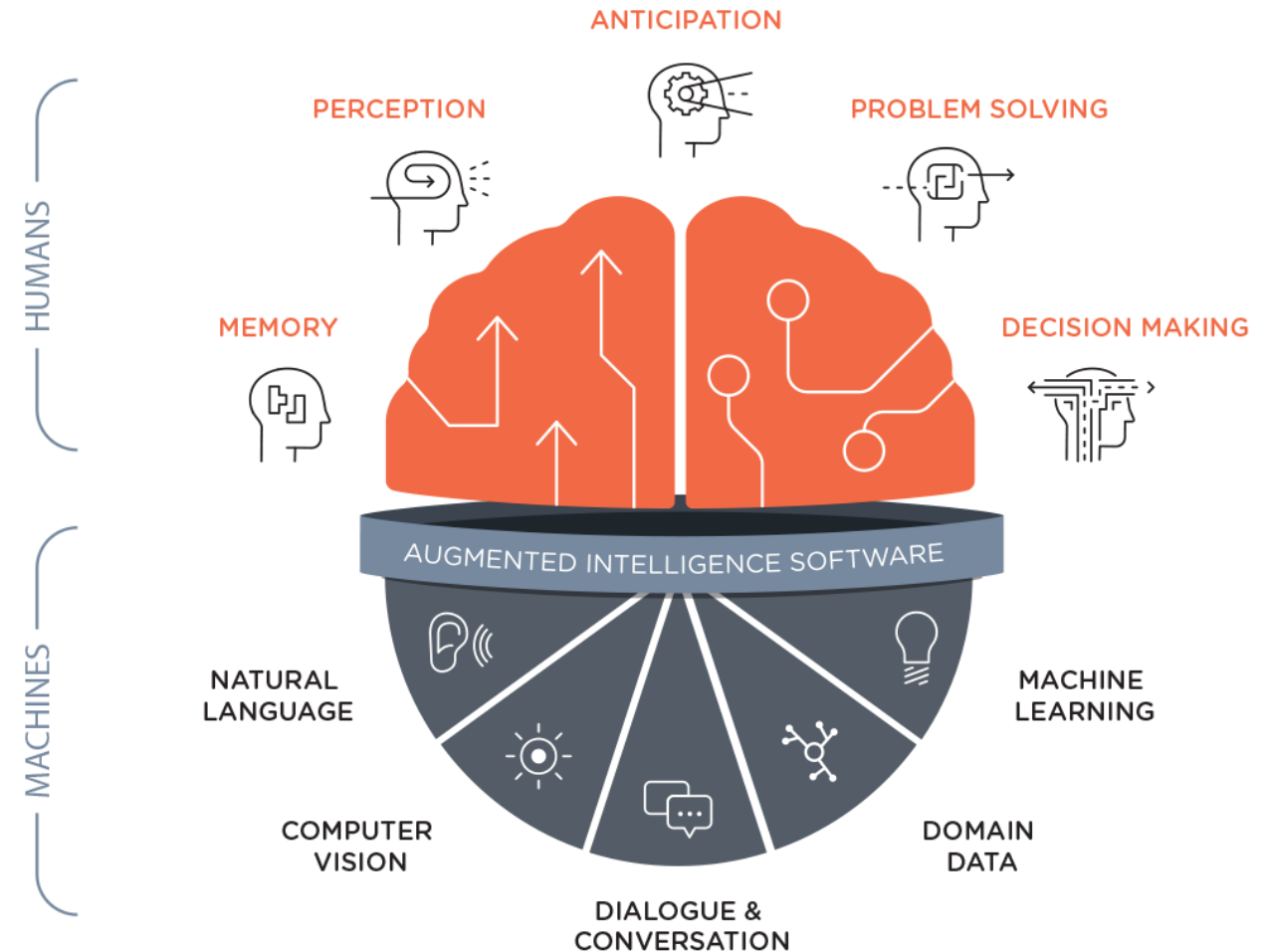
Computational Thinking

- Datalogiskt tänkande
- A **problem solving process** to describe, analyze, and solve problems such that computers can assist using techniques from computer science:
 - Give step-by-step instructions
 - Decompose problems into smaller parts
 - Find patterns
 - Create abstractions
 - Design algorithms



Why is Artificial Intelligence Different

- Scale
- Speed
- Single-mindedness
- Optimization-based
- Cannot break the rules
- No needs
- No real consequences or “skin in the game”



EU strategy for AI



A STRATEGY FOR EUROPE TO LEAD THE WAY

**Boost
technological
and industrial
capacity & AI
uptake**

**Prepare for
socio-
economic
changes**

**Ensure an
appropriate
ethical & legal
framework**



In this context: appointment of Independent High-Level Expert Group on Artificial Intelligence (AI HLEG) in June 2018

Ethics Guidelines for Trustworthy AI – Overview

Human-centric approach: AI as a means, not an end

Trustworthy AI as our foundational ambition, with three components

Lawful AI

Ethical AI

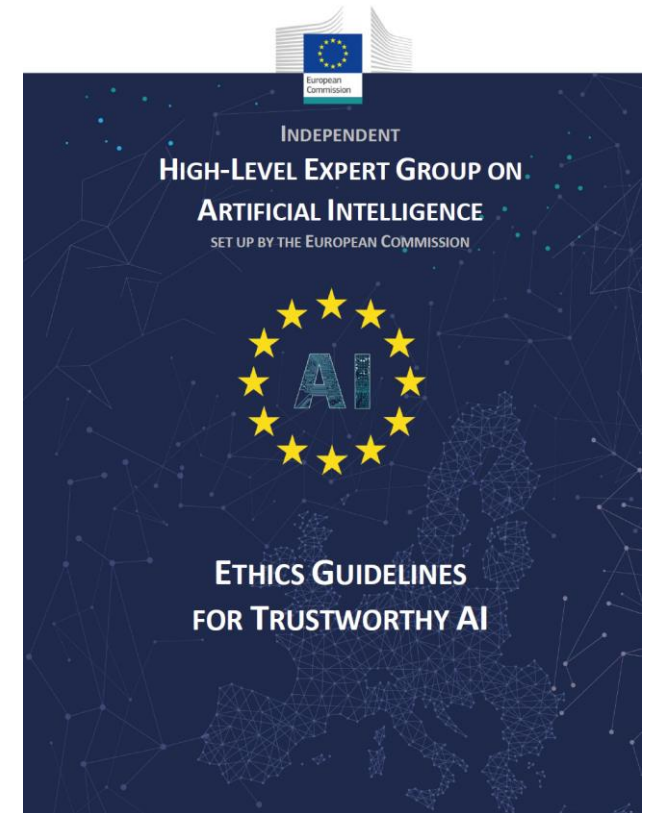
Robust AI

Three levels of abstraction

from principles
(Chapter I)

to requirements
(Chapter II)

to assessment
list (Chapter III)



Ethics Guidelines for Trustworthy AI – Principles

4 Ethical Principles based on fundamental rights



Respect for
human
autonomy

Augment, complement
and empower humans



Prevention of
harm

Safe and secure.
Protect physical and
mental integrity.



Fairness

Equal and just
distribution of
benefits and costs.



Explicability

Transparent, open
with capabilities and
purposes, explanations

Ethics Guidelines for Trustworthy AI – Requirements



Human agency and oversight



Diversity, non-discrimination and fairness



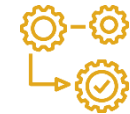
Technical Robustness and safety



Societal & environmental well-being



Privacy and data governance



Accountability

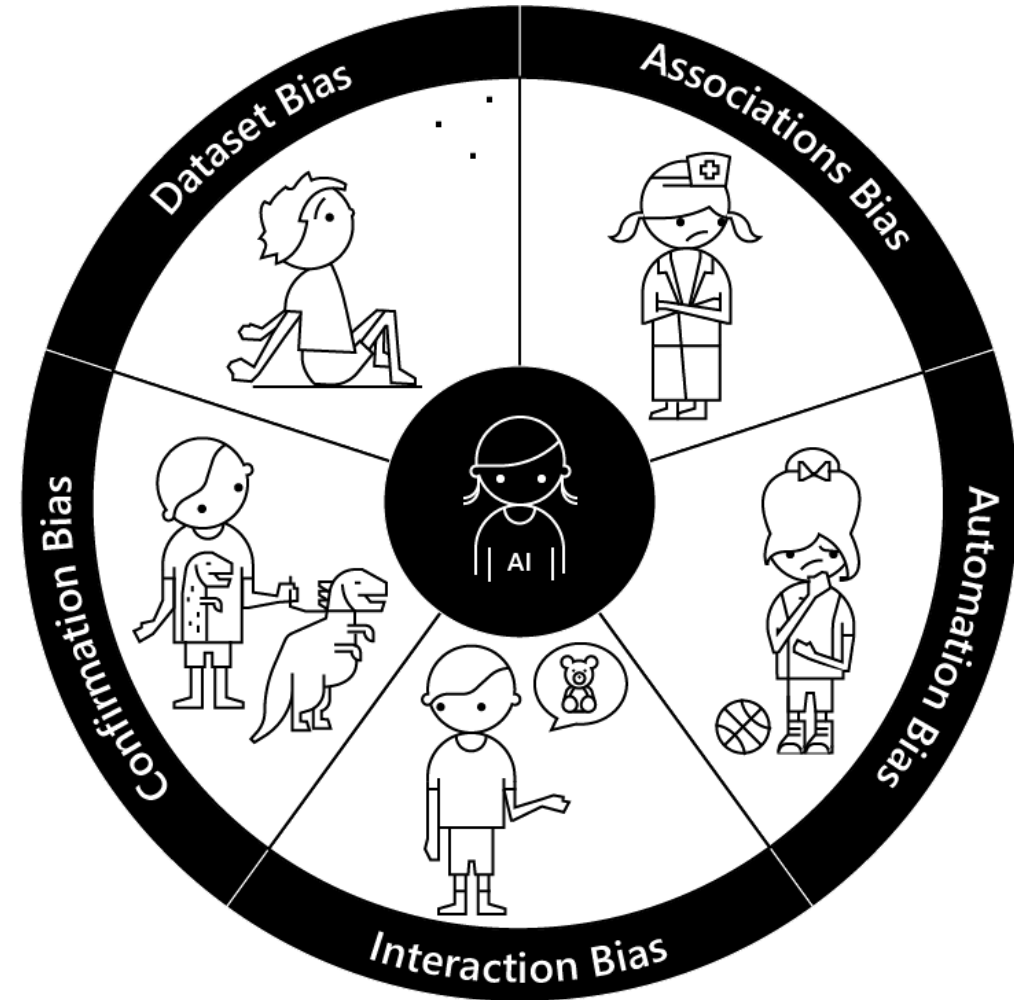


Transparency

To be continuously implemented & evaluated throughout AI system's life cycle

Bias

- **Dataset bias** – When the data used to train machine learning models doesn't represent the diversity of the customer base.
- **Association bias** – When the data used to train a model reinforces and multiplies a cultural bias.
- **Automation bias** – When automated decisions override social and cultural considerations.
- **Interaction bias** – When humans tamper with AI and create biased results.
- **Confirmation bias** – When oversimplified personalization makes biased assumptions for a group or an individual.



Research Challenges



Respect for
human autonomy

- Human-AI interaction
- Meaningful human control



Prevention of
harm

- AI Safety
- Value alignment



Fairness

- Fairness

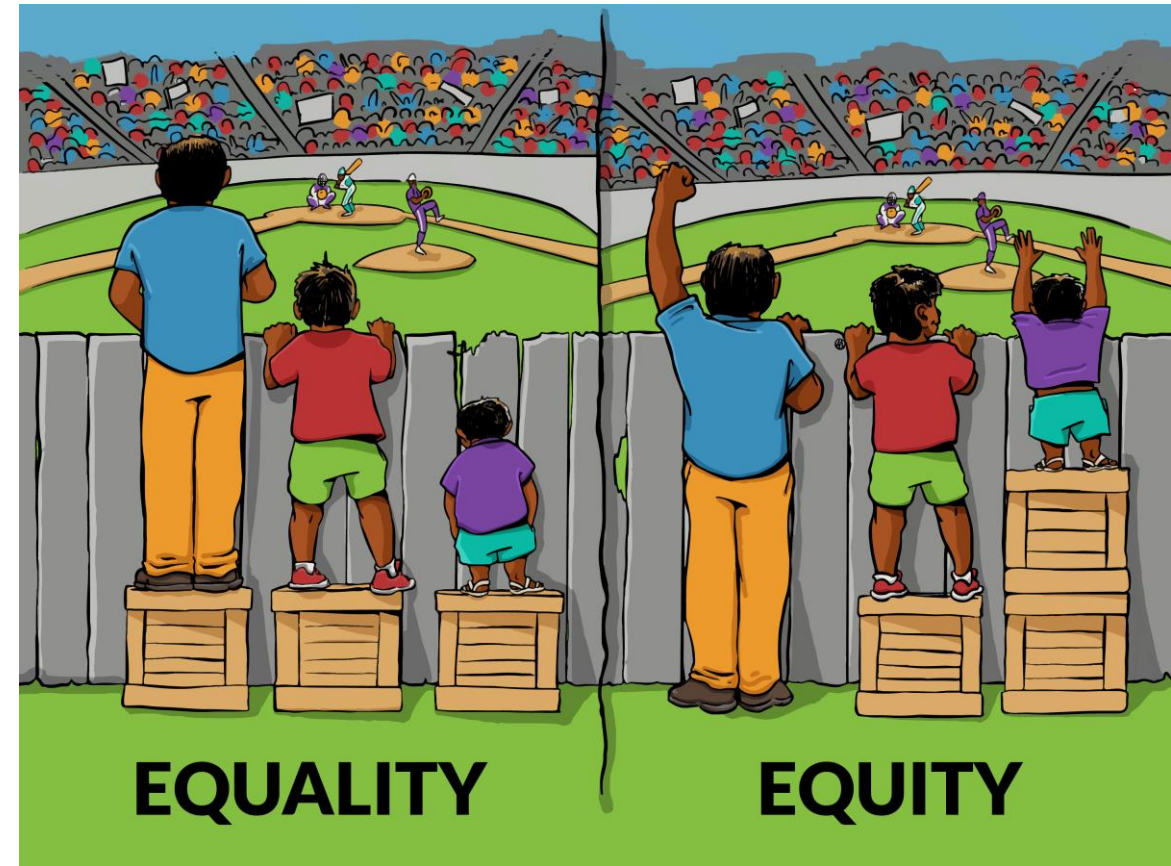


Explicability

- Accountability
- Transparency
- Explainability

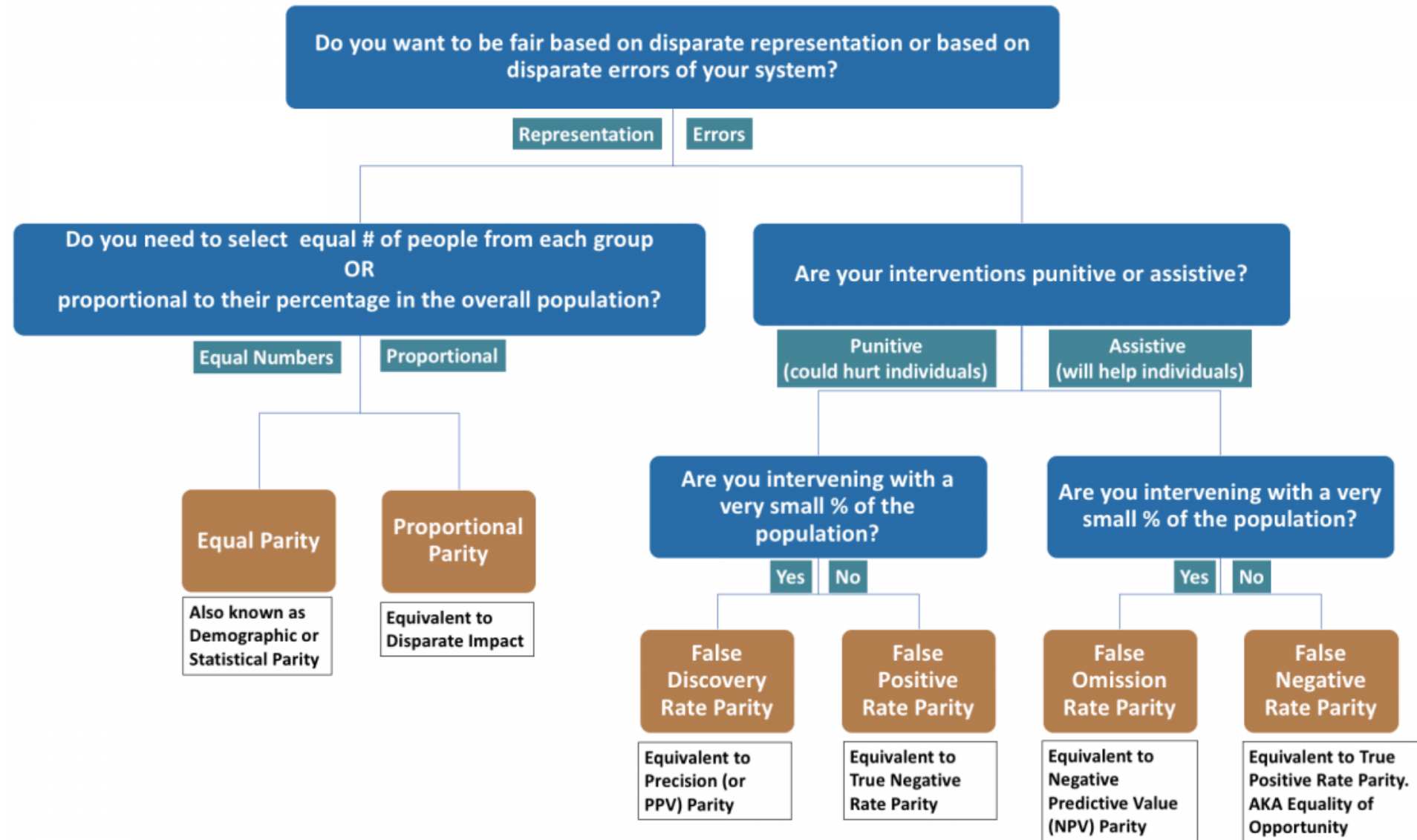
What is Fairness and Fairness-Aware ML?

- Impartial and just treatment or behavior without favoritism or discrimination.
- *Fairness-aware machine learning* algorithms seek to provide methods under which the predicted outcome of a classifier operating on data about people is fair or non-discriminatory.

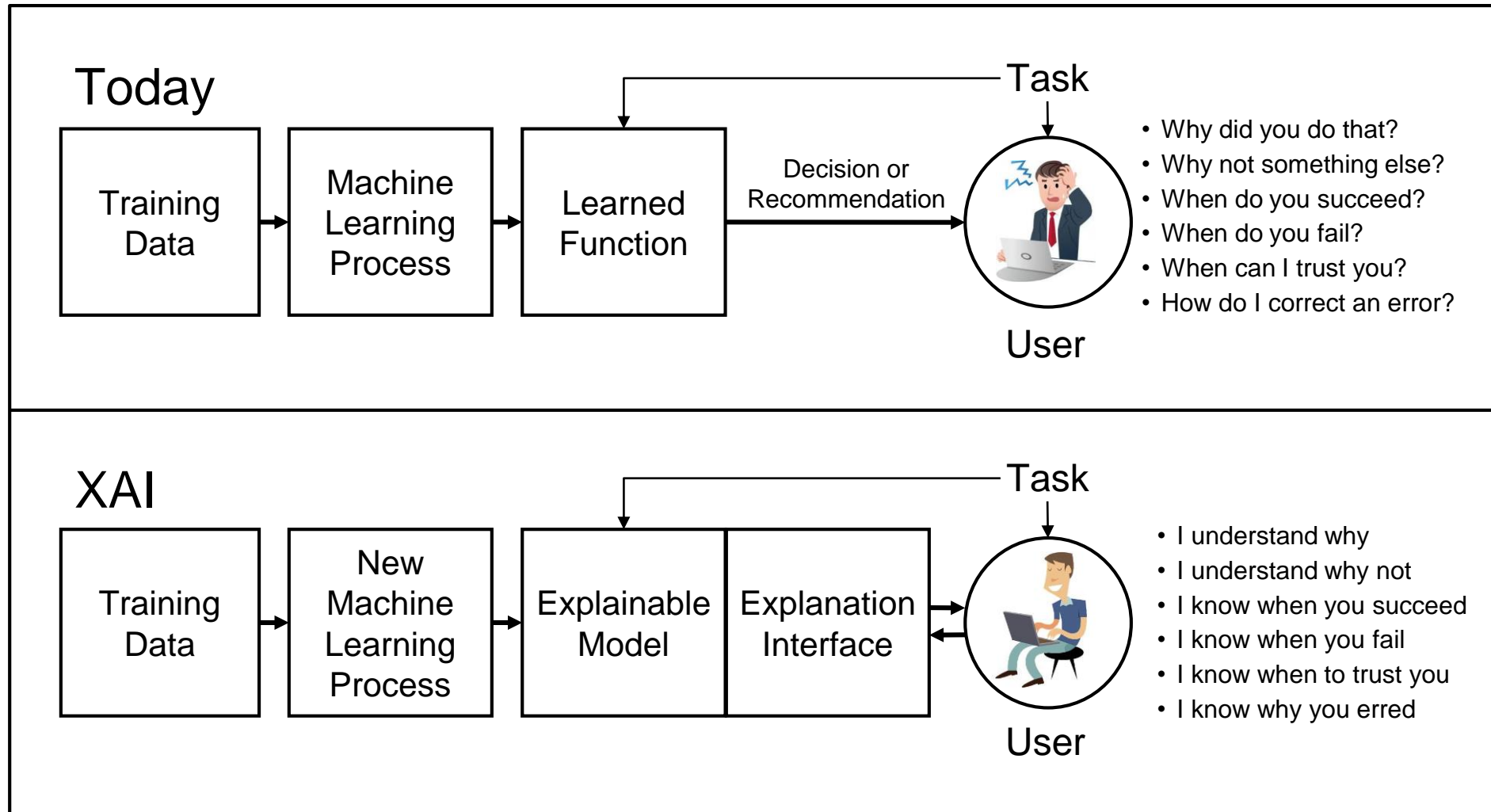


Interaction Institute for Social Change | Artist: Angus Maguire

FAIRNESS TREE



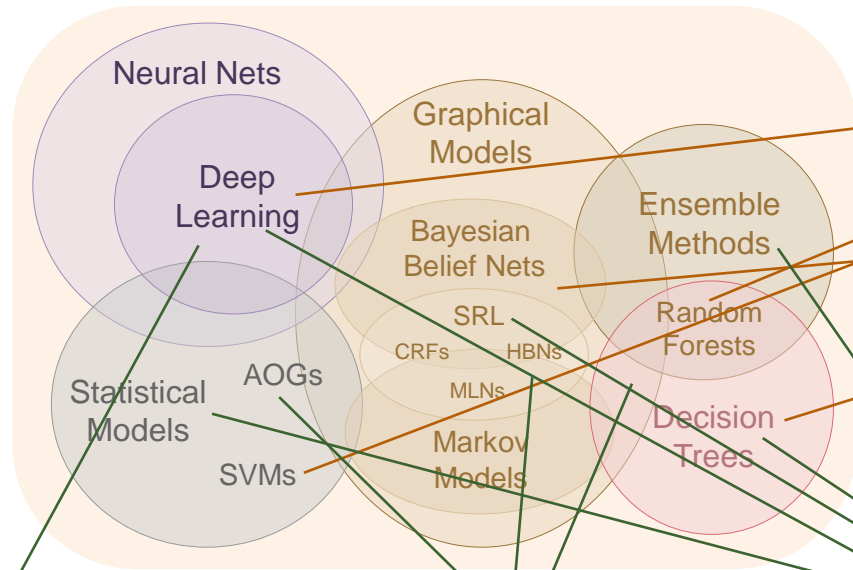
Explainable AI – The DARPA View



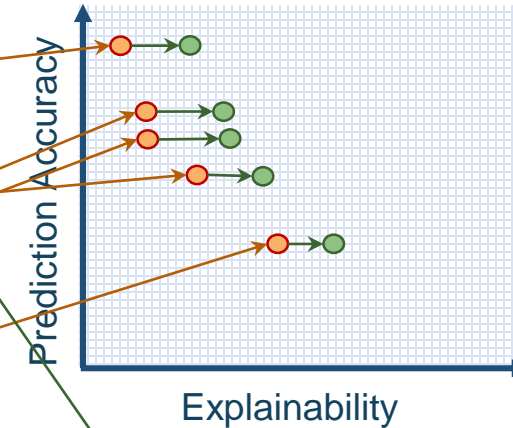
New Approach

Create a suite of machine learning techniques that produce more explainable models, while maintaining a high level of learning performance

Learning Techniques (today)



Explainability (notional)



Deep Explanation
Modified deep learning techniques to learn explainable features

The diagram shows a deep neural network with input units, hidden units, and output units. The input units are labeled 'Whiskers' and 'Claws', and the output units are labeled 'Fur' and 'Claws'. The hidden units are labeled 'A1' and 'A2'.

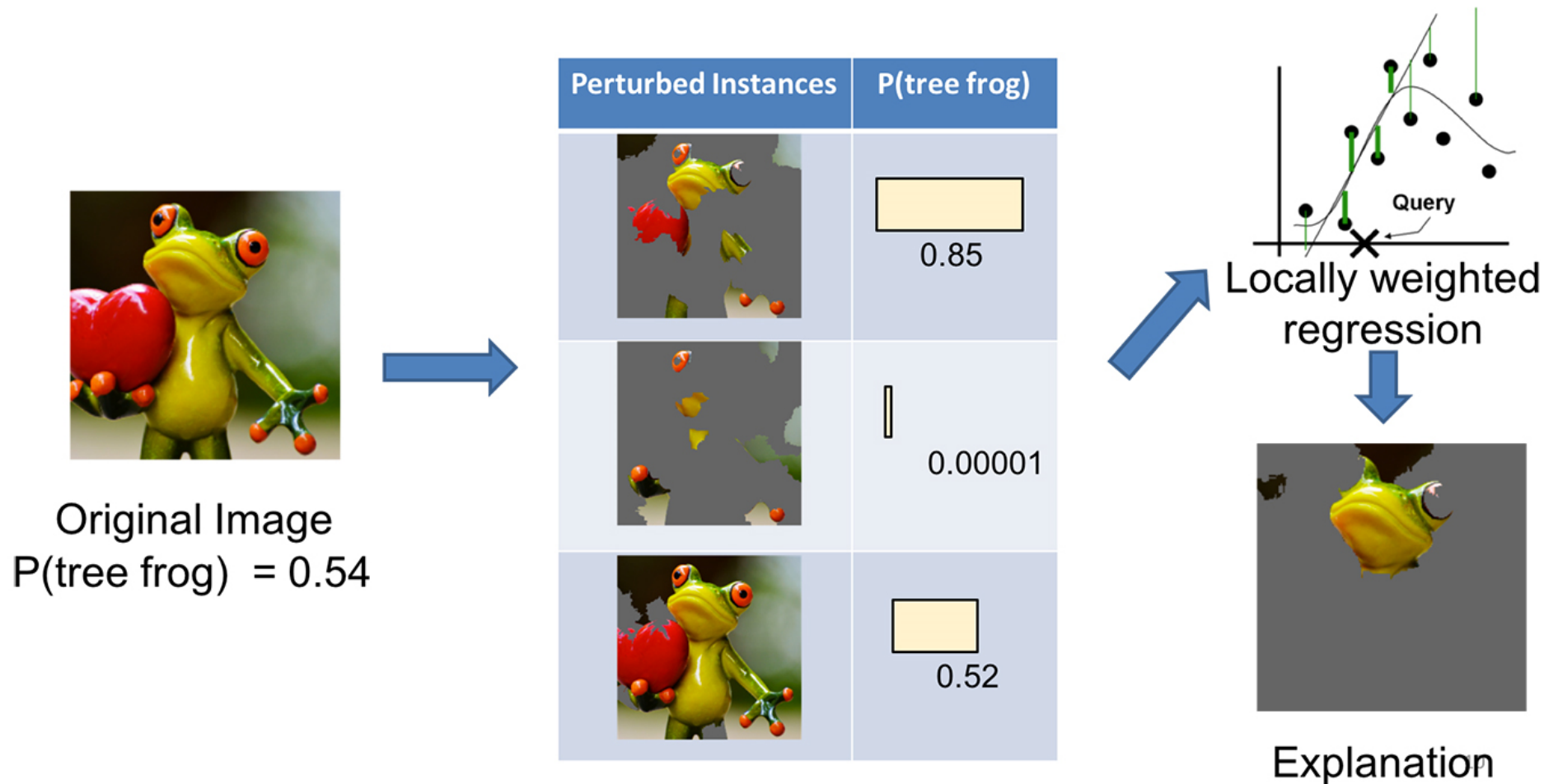
Interpretable Models
Techniques to learn more structured, interpretable, causal models

The diagram shows a decision tree model with nodes and branches. The root node is labeled 'A1'. The branches are labeled with numerical values, and the leaf nodes are labeled with numerical values.

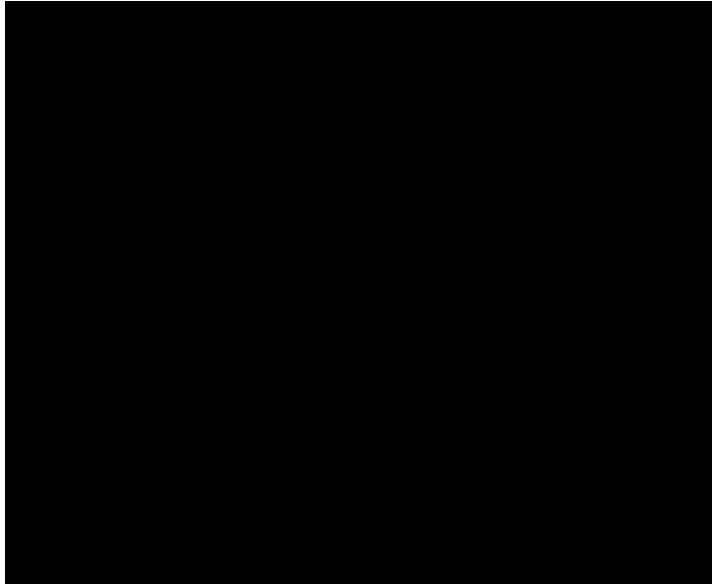
Model Induction
Techniques to infer an explainable model from any model as a black box

The diagram shows a process where a 'Model' (represented by a black box with a question mark) is used in an 'Experiment' to infer an explainable model.

LIME (Local Interpretable Model-agnostic Explanations)

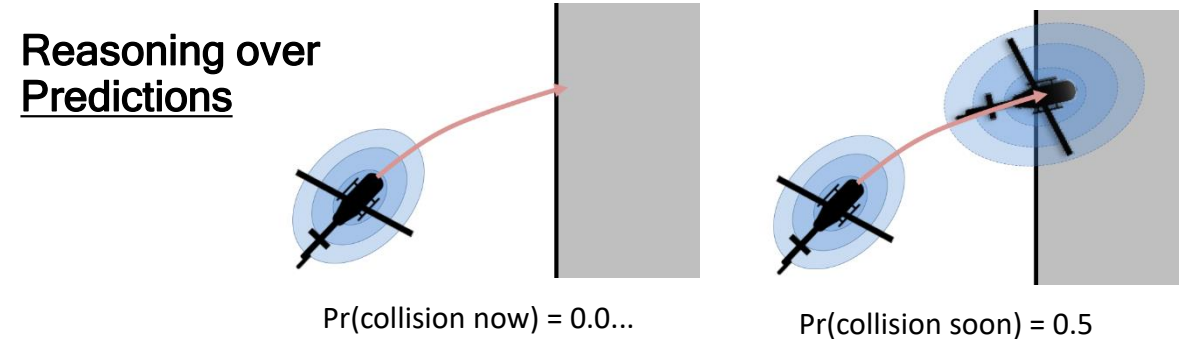
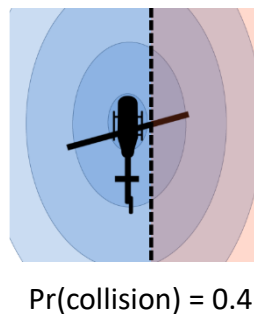
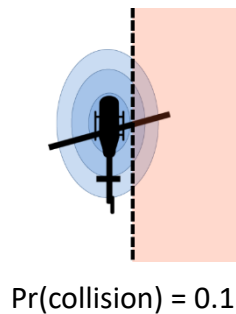
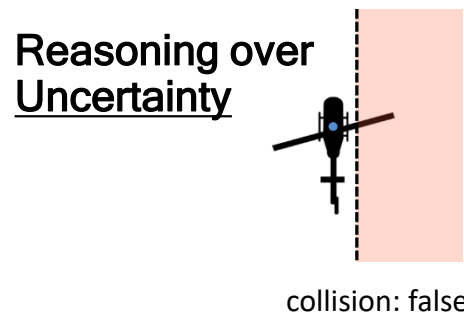


Safe Autonomous Systems / AI



If things can go wrong they probably will!

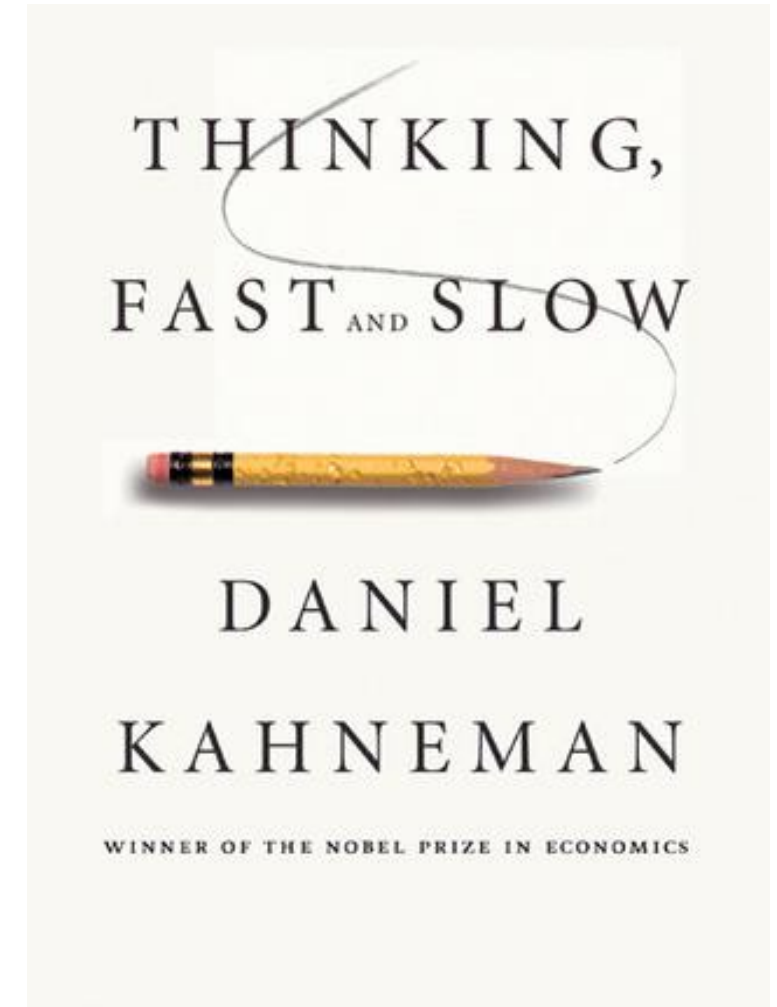
This implies the need for continual monitoring of an autonomous system and its environment in a principled, contextual, task specific manner which can be specified by the system itself!



Human and Computational Thinking

Figure 1: A Comparison of System 1 and System 2 Thinking

System 1 "Fast"	System 2 "Slow"
DEFINING CHARACTERISTICS Unconscious Effortless Automatic	DEFINING CHARACTERISTICS Deliberate and conscious Effortful Controlled mental process
WITHOUT self-awareness or control "What you see is all there is."	WITH self-awareness or control Logical and skeptical
ROLE Assesses the situation Delivers updates	ROLE Seeks new/missing information Makes decisions





Pure Logic

Pure Learning

- Slow thinking: deliberative, cognitive, model-based, extrapolation
- Amazing achievements until this day
- “*Pure logic is brittle*”
noise, uncertainty, incomplete knowledge, ...





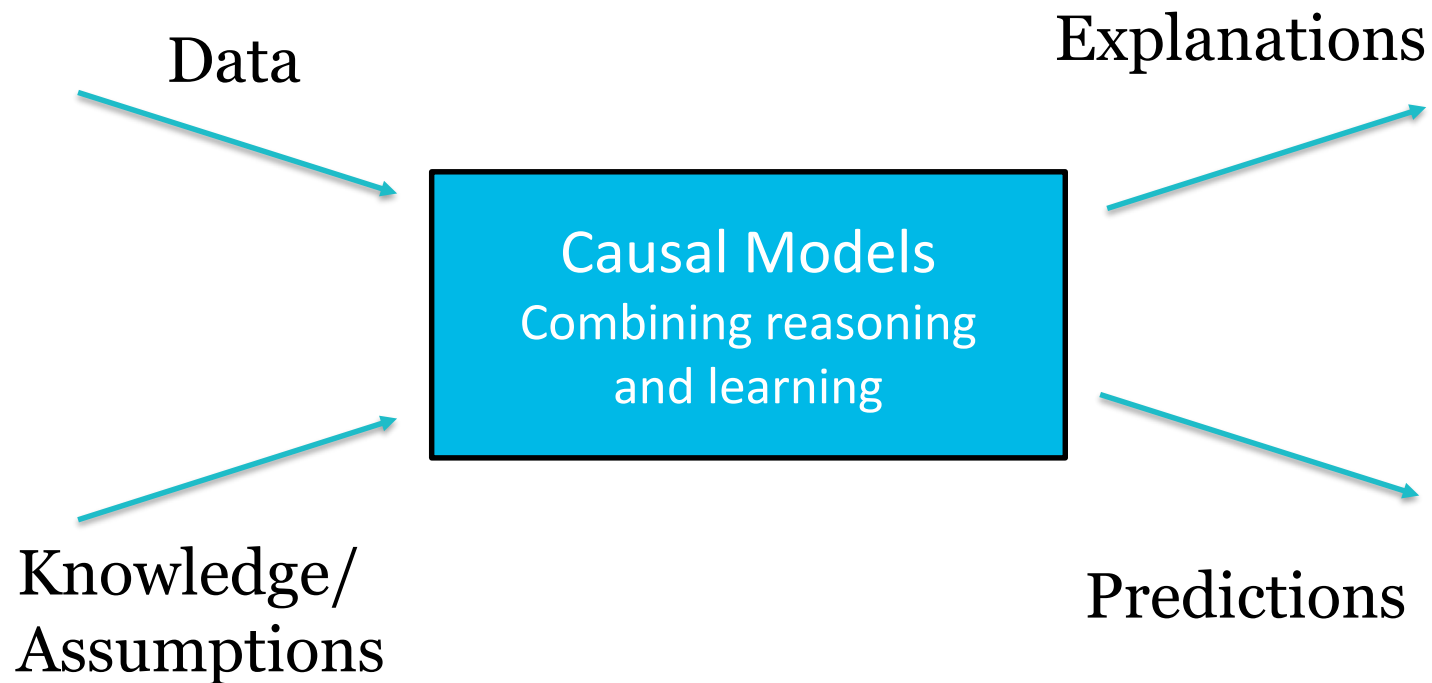
Pure Logic

Pure Learning

- Fast thinking: instinctive, perceptive, model-free, interpolation
- Amazing achievements recently
- “*Pure learning is brittle*”
 - bias, algorithmic fairness, interpretability, explainability, adversarial attacks, unknown unknowns, calibration, verification, missing features, missing labels, data efficiency, shift in distribution, general robustness and safety
 - fails to incorporate a sensible model of the world



The Way Forward



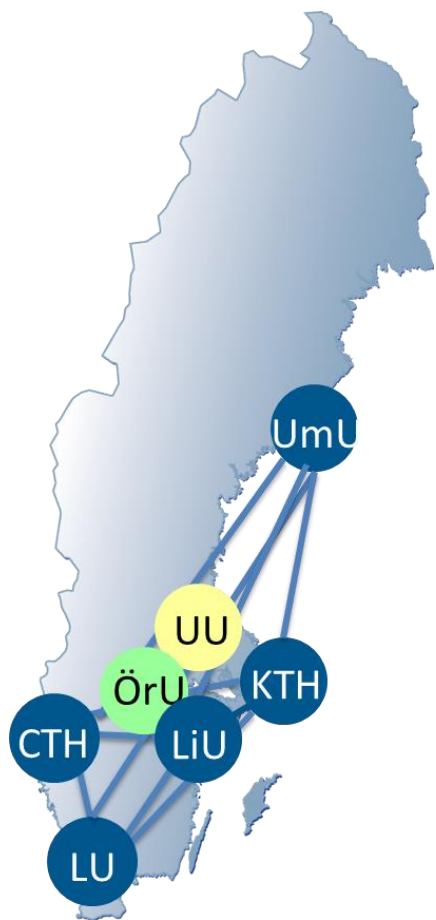


"AI systems should help
empower society,
combining the best of
technology with the best
of humanity"

<http://www.aisustainability.org/>



AI Innovation, Competence and Research in Sweden



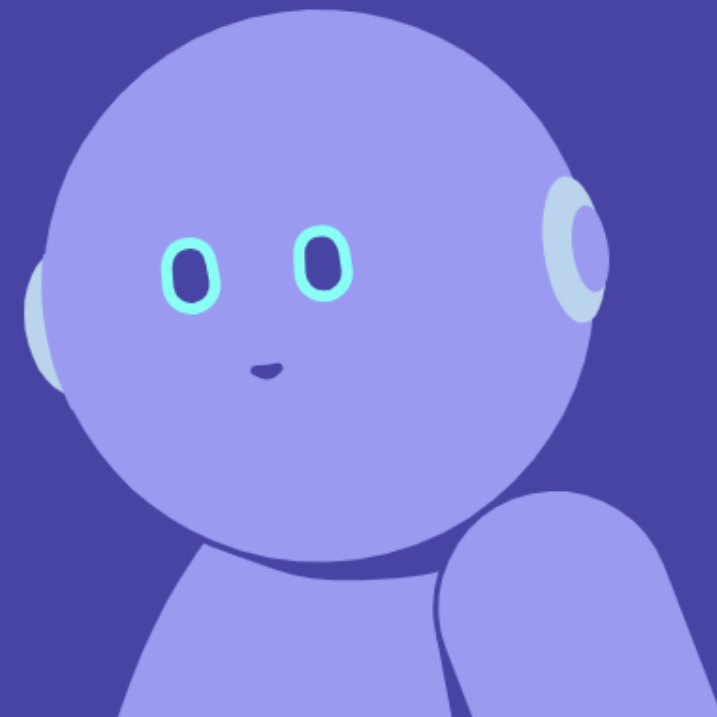
Elements of AI

Welcome to the Elements of Artificial Intelligence free online course

English ▼

Start the course

Distance course at
Linköping University
to get 2ECTS



UNIVERSITY OF HELSINKI

Reaktor

li.u LINKÖPINGS
UNIVERSITET

AI INNOVATION of Sweden



AI COMPETENCE
FOR SWEDEN

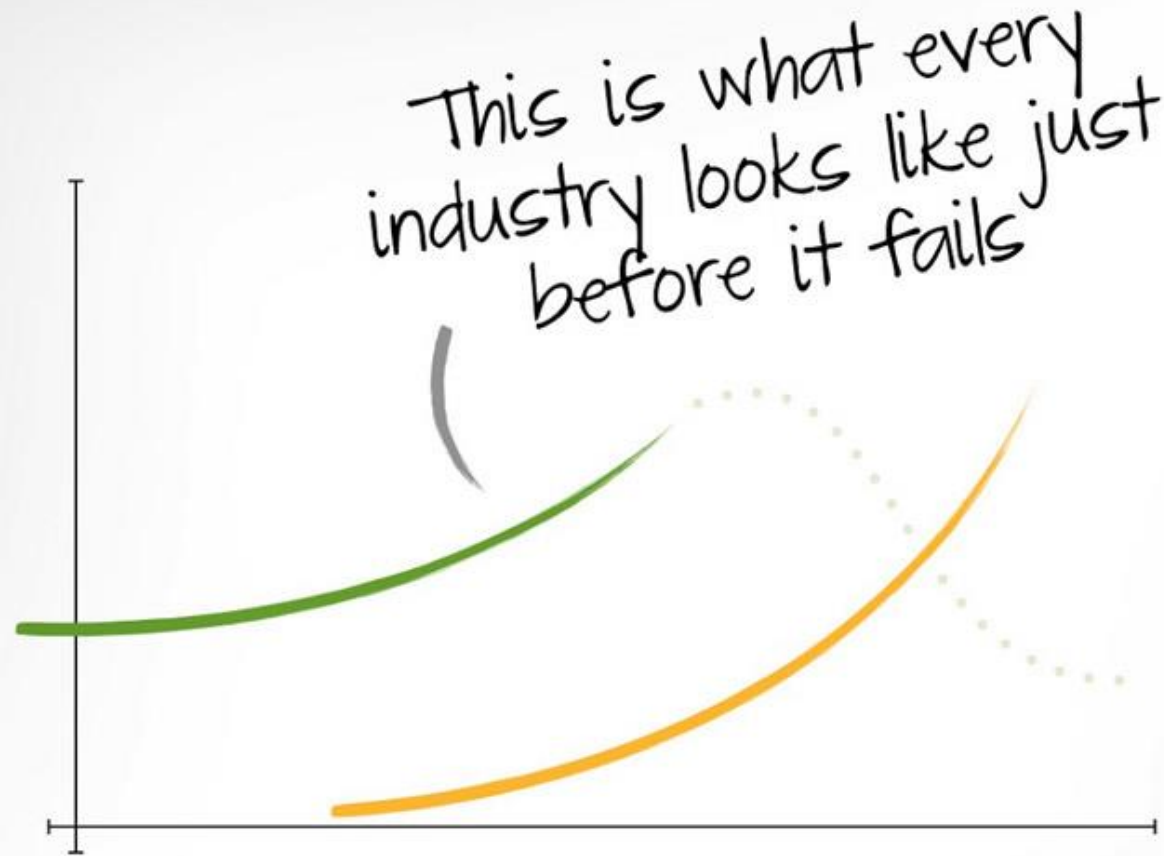
<https://www.elementsofai.se/>

Swedish launch funded by

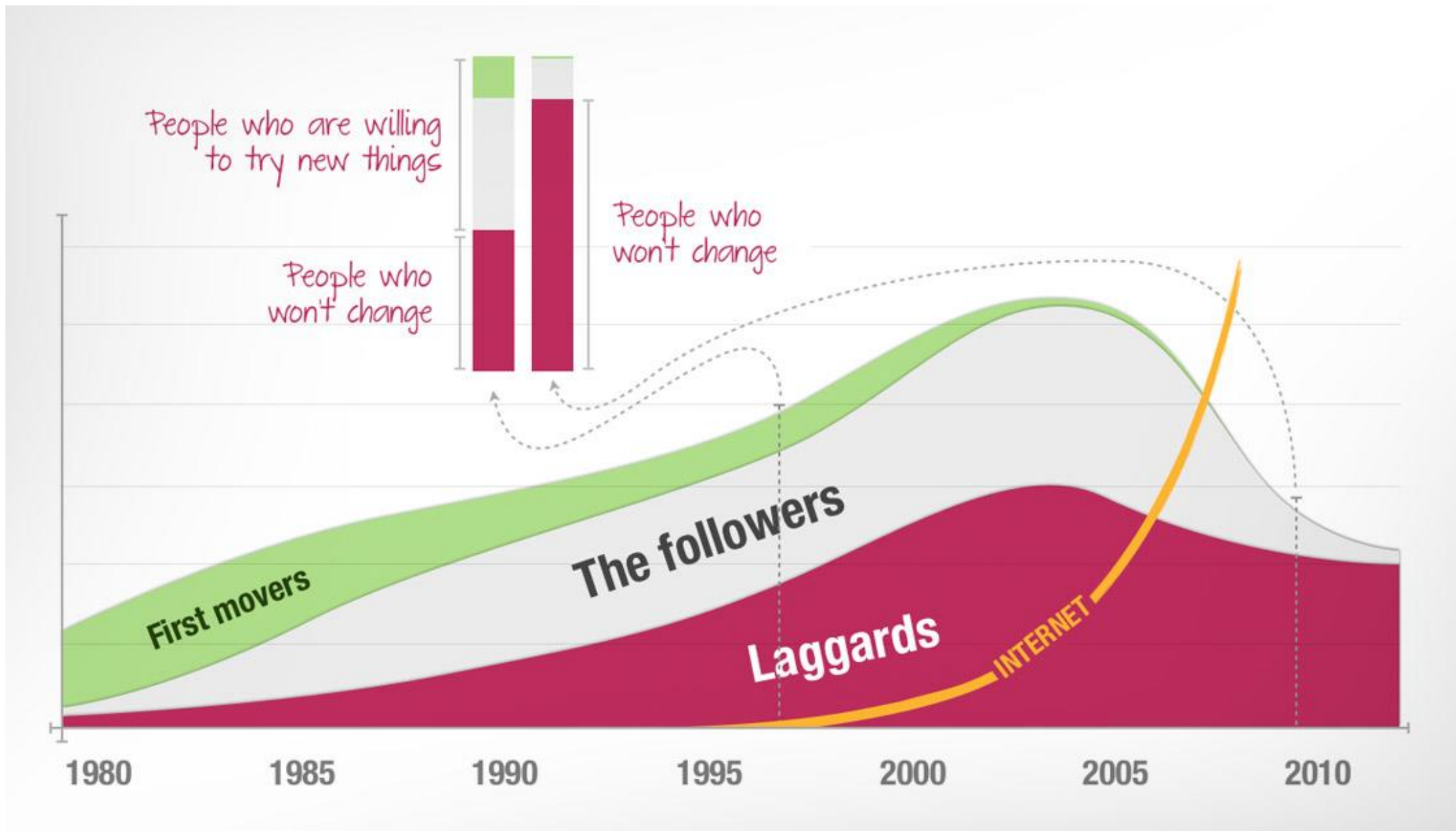
VINNOVA
Sweden's Innovation Agency

How to learn more about AI?

- Watch some overview presentations and read some popular science pieces
- Take an overview online course, e.g. <https://www.elementsofai.com/>
- Take a more technical online course and do some programming/problem solving, e.g. <https://developers.google.com/machine-learning/crash-course/>
- Take a real problem from your organization and try to solve it using AI-techniques
- Take a technical course, either online or at a university
- Build real solutions



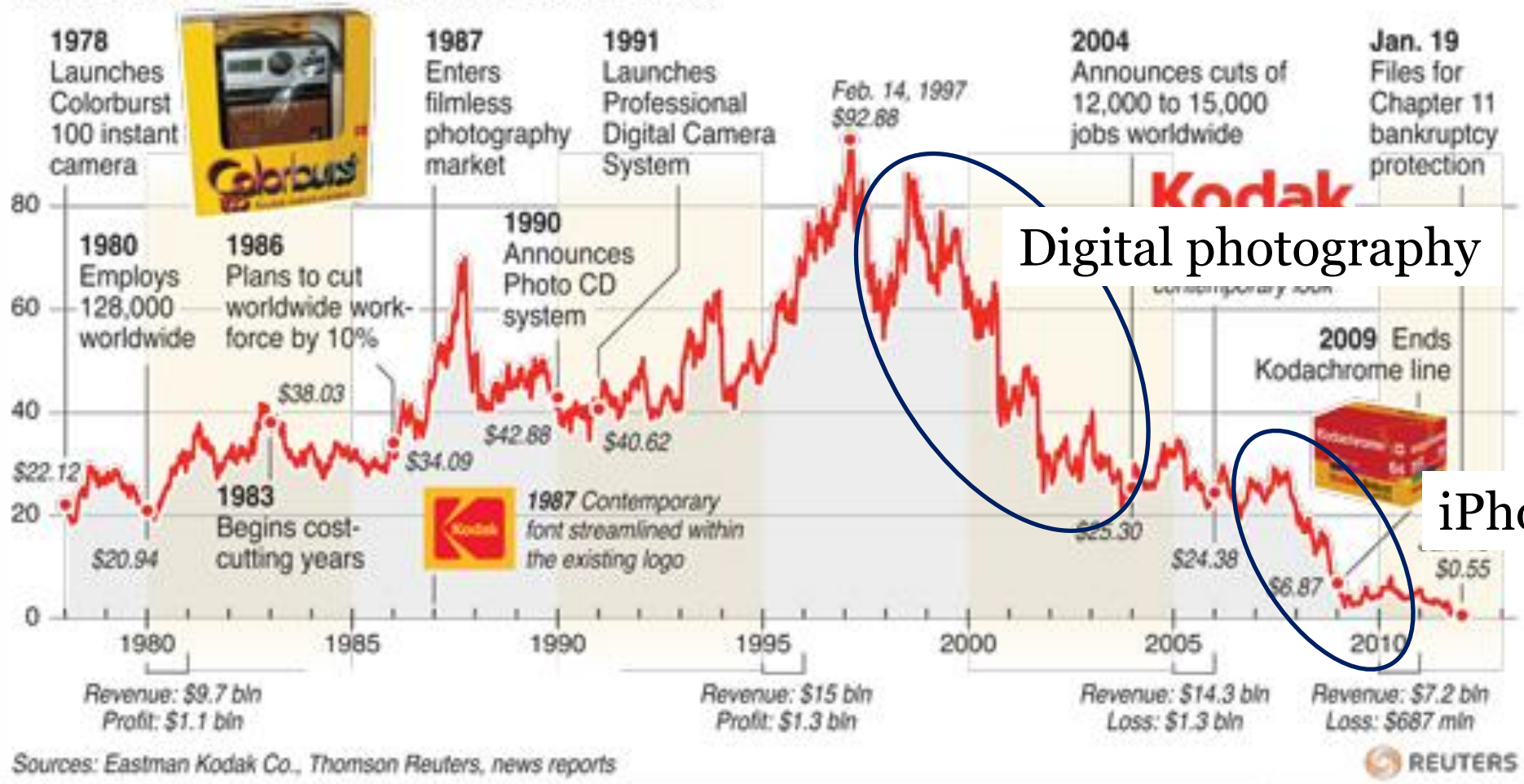
The new trends are always forming while the old world is still growing.



KODAK FILES FOR BANKRUPTCY

Eastman Kodak Co, a 130-year-old photographic film pioneer, has filed for bankruptcy protection. It said it had also obtained a \$950 million, 18-month credit facility from Citigroup to keep it going

SHARE PRICE HISTORY — WEEKLY CLOSE IN US\$



Take Away Message

- AI is about understanding intelligence and develop systems that exhibit intelligent behavior. AI will affect all aspects of our society.
- **Trust is essential.**
- To be **trustworthy** an **AI-system** should be **legal, ethical** and **robust**.
- Education and life long learning will be absolutely necessary.
 - **Take the online course Elements of AI (<http://elementsofai.se>)**
- Recommendations
 - **Learn** more! Experiment!
 - **Do** concrete projects on important topics.
 - **Scale up.**
- ***Digital tools will only provide value when you learn how to use them effectively and adapt your organization to leverage them!***
- **Human + AI**



Respect for
human autonomy



Prevention of
harm



Fairness



Explicability

Trustworthy Human-Centric AI

Fredrik Heintz, Dept. of Computer Science
Linköpings universitet
fredrik.heintz@liu.se
[@FredrikHeintz](#)